

Fraternal Twins: Unifying Attacks on Machine Learning and Digital Watermarking

Erwin Quiring, Daniel Arp and Konrad Rieck

Technische Universität Braunschweig
Brunswick, Germany

Abstract

Machine learning is increasingly used in security-critical applications, such as autonomous driving, face recognition and malware detection. Most learning methods, however, have not been designed with security in mind and thus are vulnerable to different types of attacks. This problem has motivated the research field of *adversarial machine learning* that is concerned with attacking and defending learning methods. Concurrently, a different line of research has tackled a very similar problem: In *digital watermarking* information are embedded in a signal in the presence of an adversary. As a consequence, this research field has also extensively studied techniques for attacking and defending watermarking methods.

The two research communities have worked in parallel so far, unnoticeably developing similar attack and defense strategies. This paper is a first effort to bring these communities together. To this end, we present a unified notation of black-box attacks against machine learning and watermarking that reveals the similarity of both settings. To demonstrate the efficacy of this unified view, we apply concepts from watermarking to machine learning and vice versa. We show that countermeasures from watermarking can mitigate recent model-extraction attacks and, similarly, that techniques for hardening machine learning can fend off oracle attacks against watermarks. Our work provides a conceptual link between two research fields and thereby opens novel directions for improving the security of both, machine learning and digital watermarking.

1 Introduction

In the last years, machine learning has become the tool of choice in many areas of engineering. Learning methods are thus not only applied in classic settings, such as speech and handwriting recognition, but increasingly operate at the core of security-critical applications. For example, self-driving cars make use of deep learning for

recognizing objects and street signs [e.g., 32, 66]. Similarly, systems for surveillance and access control often build on machine learning methods for identifying faces and persons [e.g. 49, 54]. Finally, several detection systems for malicious software integrate learning methods for analyzing data more effectively [e.g., 29, 30, 33].

Machine learning, however, has originally not been designed with security in mind. Many learning methods suffer from vulnerabilities that enable an adversary to thwart their successful application—either during the training or prediction phase. This problem has motivated the research field of *adversarial machine learning* which is concerned with the theory and practice of learning in an adversarial environment [27, 35, 42]. As part of this research, several attacks and defenses have been proposed, for example, for poisoning support vector machines [9, 10], crafting adversarial samples against neural networks [40, 41, 43] or stealing models from decision trees [56].

Concurrently to adversarial machine learning, a different line of research has faced very similar problems: In *digital watermarking* information is embedded in a signal, such as an image, in the presence of an adversary [15, 46]. This adversary seeks to extract or remove the information from the signal, thereby reversing the watermarking process and obtaining an unmarked copy of the signal, for example, for illegally distributing copyrighted content. As a consequence, methods for digital watermarking naturally operate in an adversarial environment and several types of attacks and defenses have been proposed for watermarking methods, such as sensitivity and oracle attacks [e.g., 1, 12, 14, 21].

Unfortunately, the two research communities have worked in parallel so far and unnoticeably developed similar attack and defense strategies. To illustrate this similarity, let us consider the simplified attacks shown in Figure 1: The middle plot corresponds to an *evasion attack* against a learning method, similar to the attacks proposed by Papernot et al. [40, 41]. A few pixels of the target image have been carefully manipulated, such that

the digit 5 is misclassified as 8. By contrast, the right plot shows an *oracle attack* against a watermarking method, similar to the attacks developed by Westfeld [62] and Cox & Linnartz [14]. Again, a few pixels have been changed; this time however to render the watermark unreadable in the target image.

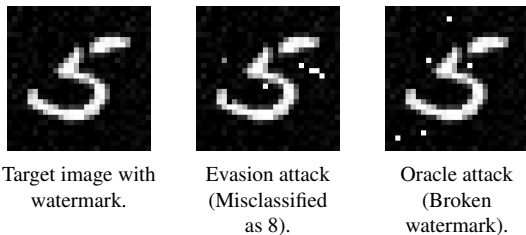


Figure 1: Examples of attacks against machine learning and digital watermarking. Middle: the target is modified, such that it is misclassified as 8. Right: the target is modified, such that the watermark is destroyed.

While both attacks address different goals, the underlying attack strategy is surprisingly similar. In fact, both attacks aim at minimally modifying the target, such that a decision boundary is crossed. In the case of machine learning, this boundary separates different classes, such as the digits. In the case of digital watermarking, the boundary discriminates watermarked from unmarked signals. Although this example illustrates only a single attack type, it becomes apparent that there is a conceptual similarity between learning and watermarking attacks.

In this paper, we strive for bringing these two research fields together and systematically study the similarities of learning and watermarking methods under an adversary’s presence with black-box access. To this end, we introduce a unified notation for attacks against learning and watermarking methods, which enables us to reason about their inner workings and abstract from the concrete attack setting. This unified view allows for transferring concepts from machine learning to digital watermarking and vice versa. As a result, we are able to apply defenses developed for watermarks to learning methods as well as transferring machine learning defenses to digital watermarking.

We empirically demonstrate the efficacy of this unified view in two case studies. First, we show that stateful defenses from digital watermarking can effectively mitigate model-extraction attacks against decision trees [56]. Second, we show that techniques for hardening machine learning with classifier diversity [6] can be successfully applied to block oracle attacks against watermarks. In addition, we provide further examples of attacks and defenses, transferable between the research fields. By doing so, we establish several links between the two research fields and identify novel directions for improving the security of both, machine learning and digital watermarking.

In summary, we make the following major contributions in this paper:

- *Machine learning meets digital watermarking.* We present a novel view on black-box attacks against learning and watermarking methods that exposes previously unknown similarities between both research fields.
- *Transfer of attacks and defenses.* Our unified view enables transferring concepts from machine learning to digital watermarking and vice versa, giving rise to novel attacks and defenses.
- *Case studies with two novel defenses.* We present and evaluate two novel defenses that are derived from our unified view and mitigate model-extraction attacks and oracle attacks, respectively.

The rest of this paper is organized as follows: In Section 2 we review the background of adversarial machine learning and digital watermarking. We introduce our unified view on both research fields in Section 3 and present case studies with defenses in Section 4. We discuss the implications of our work in Section 5 and conclude in Section 6.

2 Background

Whenever machine learning or digital watermarking are applied in security-critical applications, one needs to account for the presence of an attacker. This adversary may try to attack the learning/watermarking process and thereby impact the confidentiality, integrity and availability of the application. This section provides a basic introduction to the motivation and threat scenarios in *machine learning* and *digital watermarking*, before Section 3 systematizes them under a common notation. A reader familiar with one of the two fields may directly proceed to Section 3.

2.1 Adversarial Machine Learning

Machine learning has become an integral part of many applications in computer science and engineering, ranging from handwriting recognition to autonomous driving. The success of machine learning methods is rooted in its capability to automatically infer patterns and relations from large amounts of data [see 17, 26]. However, this inference is usually not robust against attacks and thus may be disrupted or deceived by an adversary. These attacks can be roughly categorized into three classes: *poisoning*, *evasion* and *model extraction* [42]. The latter two are the focus of our work, as they have concrete counterparts in the area of digital watermarking.

Evasion attacks. In this attack setting, the adversary attempts to thwart the prediction of a trained classifier and evade a detection. To this end, the attacker carefully manipulates characteristics of the data provided to the classifier to change the predicted class. As a result, the attack impacts the *integrity* of the prediction. For example, in the case of spam filtering, the adversary may omit words from spam emails indicative for unsolicited content [36]. A common variant of this attack type are *mimicry attacks*, in which the adversary mimics characteristics of a particular class to hinder a correct prediction [19, 52]. Evasion and mimicry attacks have been successfully applied against different learning-based systems, for example in network intrusion detection [20, 52], malware detection [25, 53, 65] and face recognition [50].

Depending on the adversary’s knowledge about the classifier, evasion attacks can be conducted in a *black-box* or *white-box* setting. In the black-box setting, no information about the learning method and its training data are available and the adversary needs to guide her attack along the predicted classes of the classifier [35, 40, 60]. With increasing knowledge of the method and data, the probability of a successful evasion rises [7]. In such a white-box setting, the adversary may exploit leaked training data to build a surrogate model and then determine what feature combinations have the most effect on prediction.

Model extraction. In this attack setting, the adversary actively probes a learning method and analyzes the returned output to reconstruct the underlying learning model [35]. This attack, denoted as *model extraction* or *model stealing*, impacts the *confidentiality* of the learning model. It may allow the adversary to gain insights on the training data as well as obtain a suitable surrogate model for preparing evasion attacks.

Depending on the output, the adversary also operates in either a *black-box* or *white-box* setting. If only the predicted classes are observable, extracting the learning model is more challenging, whereas if function values are returned or learning parameters are available the adversary can more quickly approximate the learning model. As an example, the recent attacks proposed by Tramèr et al. [56] enable reconstructing learning models from different publicly available machine learning services in black-box as well as white-box settings. Moreover, model extraction poses a serious risk to the privacy of users, as the attack may enable to derive private information from the reconstructed model [51].

2.2 Digital Watermarking

Digital watermarking allows for verifying the authenticity of digital media, like images, music or videos. Digital wa-

termarks are frequently used for copyright protection and identifying illegally distributed content [5]. Technically, a watermark is attached to a medium by embedding a pattern into the signal of the medium, such that the pattern is *imperceptible* and *inseparable*. A particular challenge for this embedding is the robustness of the watermark, which should persist under common media processing, such as compression and denoising. There exist several approaches for creating robust watermarks and we refer the reader to the comprehensive overview provided by Cox et al. [15].

As an example, Figure 2 shows a simple watermarking scheme where a random pattern is added to the pixels of an image. The induced changes remain (almost) unnoticeable, yet the presence of the watermark can be detected by correlating the watermarked image with the original watermark. Appendix A illustrates this simple watermarking scheme in more detail.

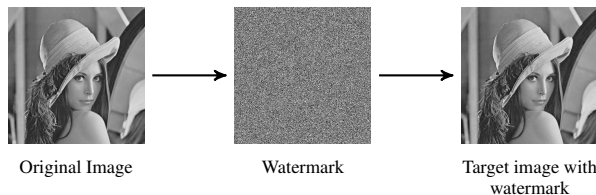


Figure 2: Example of a digital watermark. A random noise pattern is added to the image in the spatial domain. The pattern is not observable but detectable.

Similar to machine learning, watermarking methods need to account for the presence of an adversary and withstand different forms of attacks [14, 21]. While there exist several attacks based on information leaks and embedding artifacts that are unique to digital watermarking [e.g., 4, 15], we identify two attack classes that correspond to black-box evasion and model-extraction attacks.

Oracle attacks. In this attack scenario, the adversary has access to a watermark detector that can be used to check whether a given media sample contains a watermark or not [14]. Such a detector can be an online platform verifying the authenticity of images as well as a media player that implements digital rights management. Given this detector, the attacker can launch an *oracle attack* in which she iteratively modifies a watermarked medium until the watermark is undetectable. The attack thus impacts the *integrity* of the pattern embedded in the signal.

While it is trivial to destroy the pattern and the coupled signal, for example using massive changes to the medium, carefully removing the watermark while preserving the original signal is a notable challenge. As a consequence, a large variety of different attack strategies has been proposed [e.g., 12, 14, 16, 28]. A prominent example is the

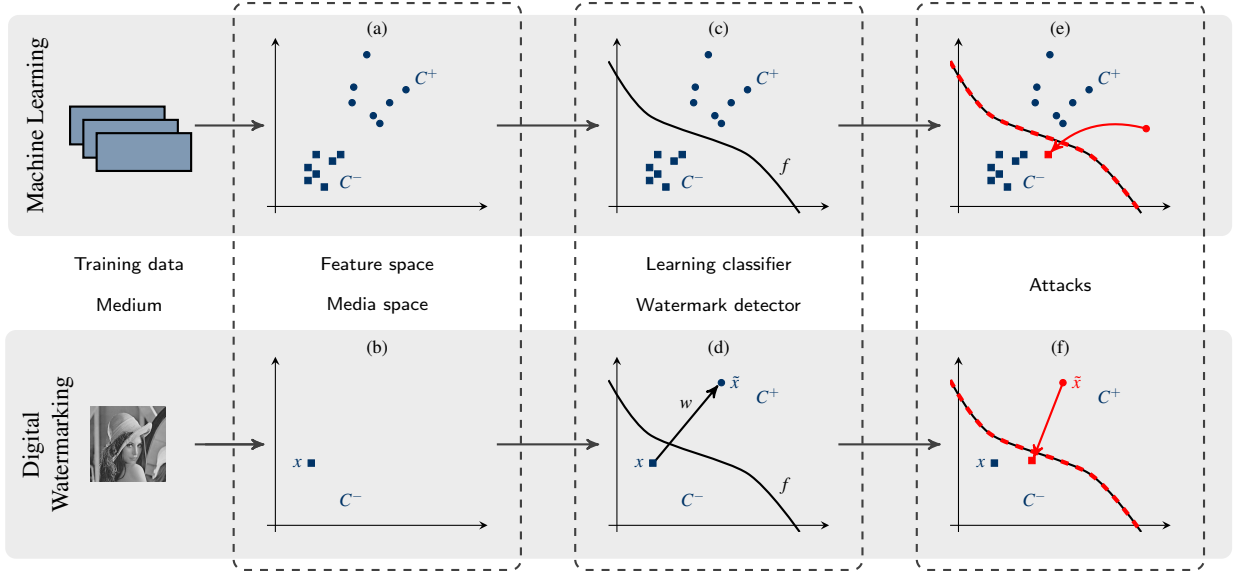


Figure 3: A unified view on machine learning and digital watermarking. Top: A machine learning setup including a feature space, a learning classifier and corresponding attacks. Bottom: A watermarking setup including the media space, the watermark detector and corresponding attacks. The red dashed line illustrates model extraction/watermark estimation, while the red arrow shows an evasion attack/oracle attack.

Blind Newton Sensitivity Attack, where no prior knowledge about the detector’s decision function is required and which has been successfully applied against several watermarking schemes (see Appendix B).

Watermark estimation. In the second setting, the adversary also has access to a watermark detector, yet her goal is not only to remove the watermark from a target medium but to estimate its pattern [11, 38]. The attack thus impacts the *confidentiality* of the watermark and not only allows to perfectly remove the pattern from the signal but also enables forging the watermark onto arbitrary other data. This *watermark estimation* poses a severe threat to watermarking methods, as it may completely undermine security mechanisms for copyright protection and access control. However, estimating the pattern embedded in a medium is difficult and requires a considerable number of queries to the watermark detector to identify discriminative features in the signal.

3 Unifying Adversarial Learning and Digital Watermarking

It is evident from the previous section that attacks against learning and watermarking methods share similarities—an observation that has surprisingly been overlooked by the two research communities [2]. Throughout this section, we systematically identify the similarities and show that it is possible to transfer knowledge about attacks and defenses from one field to the other. An overview of this

systematization is presented in Figure 3. We guide our systematization of machine learning and digital watermarking along the following four concepts:

1. *Data Representation.* Machine learning and watermarking make use of similar data representations, which enables putting corresponding learning and detection methods into the same context (Section 3.1)
2. *Problem setting.* Watermarking can be seen as a special case of a binary classification. Consequently, binary classifiers and watermarking techniques tackle a similar problem (Section 3.2).
3. *Attacks.* Due to the similar representation and problem setting, attacks overlap between both fields, as we discuss for evasion attacks (Section 3.3) and model extraction (Section 3.4).
4. *Defenses.* Defenses developed in one research field often fit the corresponding attack in the other field and thus can be transferred due to the similar data representation and problem setting (Section 3.5).

In the following, we discuss each of these concepts in more detail, where we first formalize the concept for machine learning and then proceed to digital watermarking.

3.1 Feature Space vs. Media Space

Machine learning. Learning methods typically operate on a so-called *feature space* that captures the characteristics of the data to be analyzed and learned. These features

correspond to vectors $x \in \mathbb{R}^N$ and in the case of *classification* are assigned to a class label y that needs to be learned and predicted, such as C^+ and C^- in Figure 3(a). Note that feature spaces in machine learning can also be constructed implicitly, for example using non-linear maps and kernel functions [17, 48].

Digital watermarking. Similar to machine learning, watermarking methods operate on a signal available in some underlying media space, such as the pixels of an image or the audio waves of a recording. Without loss of generality, this signal can be described as a vector $x \in \mathbb{R}^N$ and thus the media space corresponds to the feature space used in machine learning. Note that advanced watermarking schemes often map the signal to other spaces, such as frequency or random subspace domains [15, 21]. Still, the mapped signals can be described as points in a vector space.

Consequently, the feature space of a learning method is closely related to the media space used in digital watermarking. The relation remains unchanged even if a feature mapping is performed, as long as an implicit vector representation exists.

3.2 Classifier vs. Watermark Detector

Machine learning. After embedding the training data into a feature space, the actual learning process is performed using a learning method, such as a support vector machine or a neural network. In the case of classification, this learning method tries to infer functional dependencies from the training data to separate data points of different classes. These dependencies are described in a learning model w that parameterizes a decision function $f_w(x)$. Given a vector x the function $f_w(x)$ predicts a class label or a corresponding numerical prediction score.

Digital watermarking. The media space in watermarking is divided in two separate subspaces as depicted in Figure 3(d) where the marked and unmarked versions of the signal represent the two classes. Note that a robust watermark should ideally survive image processing steps, such as compression and denoising. Therefore, the watermark class implicitly contains variations as well, just as machine learning captures the variations of samples from a class through its generalization.

If we denote an unmarked signal as x and a watermarked signal as \tilde{x} , the relation between x and \tilde{x} is given by a parameter w that defines the pattern of the watermark. As a consequence, a watermark detector also employs a function $f_w(x)$ to determine which subspace a signal is in and thus whether it contains the watermark. Similar to machine learning, the function f_w may induce a linear as

well as non-linear boundary, such as a polynomial [22] or fractalized boundary [37].

Altogether, both fields perform a classification and an adversary faces the same situation: a decision boundary separates two classes either in feature or media space. Consequently black-box attacks that work through input-output observations are quite transferable between machine learning and digital watermarking. We emphasize that the boundary does not need to be the same. Our focus lies on the corresponding attack strategy. In the following sections, we discuss this similarity and provide a mapping between machine learning and watermarking attacks, which lays the ground for transferring defenses from one field to the other.

3.3 Evasion Attack vs. Oracle Attack

As the first attack mapping, we consider the pair of *evasion* and *oracle* attacks in a black-box setting. In this attack scenario, an adversary targets the integrity of the classifier’s response by inducing a misclassification from an iteratively collected set of input-output pairs. This kind of attack has been proposed for learning-based classifiers as well as watermark detectors.

Machine learning. In an evasion attack, the adversary tries to manipulate a sample with minimal changes, such that it is misclassified by the decision function f_w . Formally, the attack can thus be described as an optimization problem,

$$\arg \min_t d(t) \text{ s.t. } f_w(x+t) = y^* , \quad (1)$$

where $d(t)$ reflects the necessary changes t on the original sample x to achieve the wanted prediction y^* . Depending on the particular output of the learning classifier, the attacker can run different attack strategies:

- *Numerical output.* In this case, the classifier returns a prediction score $f_w(x)$ and the attacker tries to mislead the classifier with as minimal changes as possible. For example, the adversary can perform a gradient descent in the direction of the decision boundary to determine the features that have the most effect on the classification [7, 41].
- *Binary output.* In the second case, the classifier only returns the predicted class label. This clearly restricts the attacker’s capabilities, since not every change of a feature influences the classifier’s output. Still, an adversary can perform a line search through the binary responses to locate the boundary’s position [35]. An attacker can also learn a substitute model based on a set of queries that approximates

the original model [39, 40]. The attack iteratively sends new queries in regions where the substitute model is less confident. This allows an evasion even with highly non-linear models, such as deep neural networks.

Depending on the concrete scenario the attacker might also have to satisfy additional constraints. For instance, it might not be sufficient to just cross the decision boundary. Instead, the modified sample also needs to be located inside the distribution of the target class [7].

Digital watermarking. In an oracle attack, an adversary tries to disturb or even remove the watermark embedded in a medium. The attack setting is closely related to evasion. Formally, the underlying optimization problem is given by

$$\arg \min_t d(t) \text{ s.t. } f_w(\tilde{x} + t) = y^- , \quad (2)$$

where $d(t)$ reflects the changes t on the watermarked signal \tilde{x} and y^- corresponds to no detection. The optimization problem is identical to the one given in Eq. (1), so that the adversary can apply similar attack strategies. We can again categorize these strategies depending on the output returned by the watermark detector.

- *Numerical output.* In this case, the watermark detector outputs the score $f_w(\tilde{x})$ of the decision function. As for evasion, the adversary can perform an attack based on gradient descent to remove the watermark \tilde{x} from the image with as little changes as possible [14].
- *Binary output.* In this setting, the adversary has only access to the binary output of the watermark detector. As a watermark detector usually does not need to return more information than necessary, the watermarking literature generally focuses on this setting. Similar to the evasion case, it is possible to perform a line search to locate the decision boundary and remove the watermark [e.g., 12]. We present a novel defense against this type of attack in Section 4 which is inspired by concepts from adversarial machine learning.

Due to the equivalent objectives in Eq. (1) and (2), attack strategies from machine learning are transferable to watermarks and vice versa. Take, for instance, the state-of-the-art Blind Newton Sensitivity Attack [12] that solves the optimization problem from Eq. (2) by performing a gradient descent based on binary outputs. This makes the attack also applicable against non-linear learning classifiers. Appendix B recaps the attack procedure in more detail. The optimal solution is

guaranteed for convex boundaries, but suitable results are also reported for non-linear watermarking schemes by following the boundary’s envelope [12, 13].

We conclude that evasion attacks on classifiers and oracle attacks on watermark detectors share fundamental similarities in the black-box setting. The underlying optimization problems are identical and thus several of the existing attack and defense strategies can be directly exchanged from one research area to the other.

3.4 Model Extraction

As the second attack mapping, we consider the pair of *model extraction* and *watermark estimation*. In the black-box scenario, the adversary aims at compromising the confidentiality of a learning model or digital watermark by sending specifically crafted objects to a given classifier/detector and observing the respective binary output.

Machine learning. Model-extraction attacks center on an effective strategy for querying a classifier, such that the underlying model can be reconstructed with few queries. For instance, Tramèr et al. [56] have recently demonstrated this threat by stealing models from cloud platforms providing machine learning as a service. In contrast to evasion, the extraction of the learning model w enables the adversary to also reconstruct the decision function f_w and to apply it to arbitrary data. We can differentiate two attack strategies here:

- *Approximation.* In the first case, an attacker collects a number of input-output pairs with queries either scattered over the feature space or created adaptively [39, 40, 56]. These observations allow the adversary to learn an own surrogate model. While this strategy is easy to implement, it only yields an approximation of the original model, which becomes more accurate the more observations are conducted.
- *Reconstruction.* The localization of points on the decision boundary through a line search enables an exact reconstruction of the model. For example, an adversary can reconstruct a linear classifier by performing a line search in each feature direction from a fixed position [35]. The distance from this position to the located boundary leaks the respective feature weight in that direction. The extraction against non-linear classifiers such as decision trees also exploits localized boundary points for reconstruction [56]. We discuss the latter attack in more detail in Section 4 when presenting a novel defense against it, inspired by concepts from digital watermarking.

Defenses	Research Field	
	Adversarial Learning	Watermarking
Randomization	Random Subspace Method [8] Randomized Ensemble [8, 31] —	Random Subspace Method [18, 57] Randomized Boundary [21, 34] Union of Watermarks [21]
Complex Boundary	Non-Linearity [6, 47] Classifier Diversity [6] —	Non-Linearity [21–23, 37] — Snake Traps [21]
Stateful Analysis	— — —	Security Margin [1, 55] Line Search Detection [1] Locality-Sensitive Hashing [58]

Table 1: Comparison of defense techniques introduced by adversarial learning and digital watermarking.

Digital watermarking. Watermark estimation represents the counterpart to model extraction. In this attack scenario, the adversary seeks to reconstruct the watermark w from a marked signal \tilde{x} . If successful, the adversary is not only capable of perfectly removing the watermark w from the signal \tilde{x} , but also of embedding w in other signals, thereby effectively creating forgeries.

Similar to the model extraction case, estimation attacks in the watermarking literature are based on localizing boundary points where the signal just crosses the detector’s decision boundary [11, 37]. A watermark with a linear boundary, for instance, can be recovered from a linear number of discovered boundary points. Choubassi and Moulin present a variant of this estimation attack to find boundary points that reduce the effort of the subsequent watermark estimation [11]. This approach already comes very close to the work of Lowd and Meek [35] from adversarial machine learning.

As digital watermarking and machine learning use a decision boundary to separate inputs, a natural attack strategy consists in localizing this boundary through queries and then combining the gathered points to reconstruct the model or watermark.

3.5 Defenses

The communities of both research fields have extensively worked on developing defenses to fend off the attacks presented in the previous sections. However, it is usually much easier to create an attack that compromises a security goal, than devising a defense that effectively stops a class of attacks. As a result, several of the developed defenses only protect from very specific attacks and it is still an open question how learning methods and watermark detectors can be generally protected from the influence of an adversary. In this section, we provide an overview of current defenses and identify similarities as well as

interesting directions for transferring a defense strategy from one field to the other (see Table 1). We also include defenses from adversarial learning that were initially presented against informed attacks, but also work when an adversary acts in a black-box setting.

Randomization. A simple yet effective strategy to impede attacks against classifiers and watermark detectors builds on the introduction of randomness. Several techniques have been proposed in both fields which add elements of randomization to the learning or detection process. While these defenses cannot rule out successful attacks, the induced indeterminism obstructs simple attack strategies and requires more thorough concepts for evasion or model extraction.

In machine learning, *randomized ensemble learning* has been proposed for implementing this defense strategy [8, 44]. Each classifier in an ensemble is built with a random subset of the training data and the prediction is retrieved by aggregating the output of all classifiers. As a consequence, the adversary has to attack different classifiers at the same time [8]. Alternatively, the features selected to train each classifier can be randomized, such that an adversary cannot be sure whether a specific feature has an influence on the returned classifier output [60]. This is known as the *random subspace method* in the machine learning field. Overall, randomizing the training data and features raises the bar for all discussed attacks, since the adversary has to spend more effort into gaining background knowledge on the underlying decision boundary and model.

Similar techniques have been proposed to defeat attacks on watermarking detection. In particular, a detector can be hardened by creating a *randomized region* around the decision boundary where the detector returns arbitrary outputs [21, 34]. This misleads the inherent line search in attacks that localize the boundary in this way [11, 12]. Moreover, equivalent to the random subspace method

in machine learning, several works in the field of watermarking propose to randomly divide the image pixels into subsets and aggregate the classifier output from each subset [18, 57].

In addition, the *Broken Arrows* watermarking scheme creates several watermarks that form a *union of watermarks*. During detection, only the watermark with the smallest distance to the current signal is applied [21]. This mitigates the risk that an adversary could compare multiple images with the same watermark. This defense has not been applied to learning methods yet. It would correspond to an ensemble of classifiers where only one classifier is applied during prediction depending on the input sample.

Complex boundary. Another strategy for obstructing attacks is the selection of a complex decision boundary. Without sufficient knowledge on the structure of the boundary, it is difficult to attack a method and the adversary is required to invest more resources to circumvent the defense. However, increasing the complexity of a decision function is not trivial, as a fine-grained boundary can also lead to overfitting and further possibilities for evasion [see 41].

Recent work on adversarial machine learning thus proposes to increase the complexity of the decision boundary but at the same time to enclose the learned data tightly. In the case of malware detection, this implies that an evasion attack needs to contain plausible features of the benign class without losing the malicious functionality. Russu et al. implement this defense strategy using non-linear kernel functions [47], while Biggio et al. realize a tighter and more complex boundary through the combination of two-class and one-class models [6]. Although invented against attackers with a surrogate model, these countermeasures also tackle black-box attacks that need to probe the feature space with queries outside the training data distribution. We use this strategy in Section 4.1 to address a watermark oracle attack.

Similarly, the watermarking community has examined non-linear boundaries to defend against oracle and watermark-estimation attacks [22, 23, 37]. These boundaries range from polynomial to non-parametric fractals and obstruct e.g. attacks that estimate the decision boundary. In addition, Furon and Bas have introduced small indents called *snake traps* at the decision boundary in order to stop attacks based on random walks along the detection region [16, 21].

While all of the defenses listed in Table 1 using a complex boundary render attacks more difficult, it is often not clear whether they provide protection in the long run. For example, boundaries based on fractals and snake traps block some of the attacks presented in Section 3, yet ap-

proximations of the decision boundary are still possible and might be used for successful attacks [12].

Stateful analysis. If the learning method or watermark detector is outside of the attacker’s control, an active defense strategy becomes possible, in which the defender seeks to identify sequences of malicious queries. For instance, a cloud service providing machine learning as a service may monitor incoming queries for patterns indicative of evasion and model-extraction attacks.

While this concept has not yet been examined in adversarial machine learning, stateful analysis of queries has been successfully applied in digital watermarking for detecting oracle and watermark-estimation attacks [1, 55, 58]. These defenses exploit the fact that an adversary first needs to perform an unusual number of queries close to the boundary in order to exactly locate its position. Thus, it is possible to detect attempts to attack the detector and infer the decision boundary.

Consequently, Table 1 shows that stateful defenses have only been applied to watermarking schemes, providing the opportunity for constructing novel defenses for learning methods. We show in a case study in Section 4 that model-extraction attacks can be mitigated with the security margin concept if the learning system is not under full control of the adversary and it is possible to monitor incoming queries.

4 Transfer of Attacks and Defenses

We proceed to present two case studies that exemplify how concepts from one research field can be transferred to the other. As the first case study, we apply a concept proposed by Biggio et al. [6] for securing machine learning to a watermark detector. We demonstrate that the resulting detector mitigates a state-of-the-art oracle attack. In the second case study, we apply the concept of stateful detection to a machine learning method and show that this combination effectively tackles model-extraction attacks against decision trees. While these case studies focus on two particular defenses, we encourage the communities to work with each other and therefore summarize further directions for research in Section 5.

4.1 From Machine Learning to Watermarking

In our first case study, we consider a recent defense from machine learning that increases the complexity of the decision boundary by combining a two-class and one-class model [6]. Figure 4(a) schematically illustrates the concept of this defense, which effectively creates a blend between two independent learning methods:

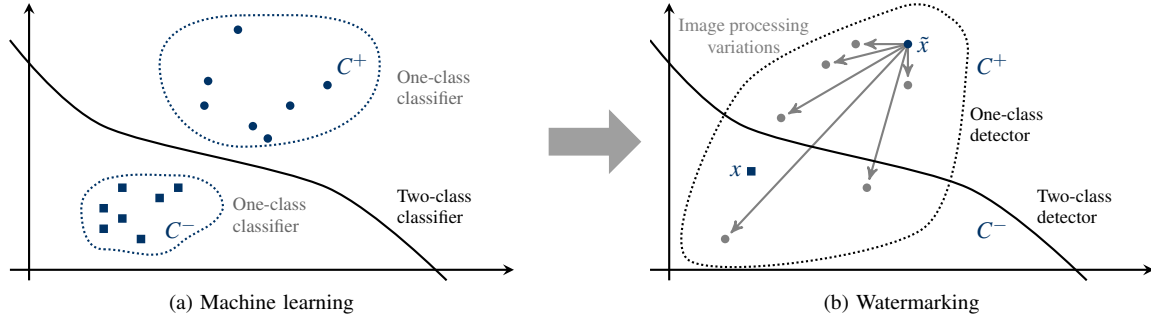


Figure 4: Transfer from machine learning to watermarking: The left plot shows the combination of a two-class and one-class learning model to build a so-called 1.5-class classifier. The right plot shows our novel defense for watermarking that also combines a two-class and one-class detector.

- *Two-class models.* The objective of this learning setting is to discriminate objects from two classes. However, unpopulated regions, such as the top-left corner in Figure 4(a), are not excluded from this classification, which leads to a weak spot: an attacker can try to evade the classification by creating arbitrary samples on the selected subspace of the decision function—irrespective of the distribution of the target class.
- *One-class models.* In this learning setting, only one class is modeled and the decision boundary separates this class from the rest of the feature space. Figure 4(a) exemplifies the concept by showing the learned boundary around the classes. One-class models enable identifying implausible points, as they tightly enclose the training data and thereby help to mitigate the weak spot of common two-class models.

If we combine the decision boundary of both models, we obtain a hybrid form denoted as “1.5-class classifier”. This classifier discriminates two classes but also requires these classes to lie within specific regions of the feature space. As a result, evasion attacks become more difficult, since an adversary needs to stay within the one-class regions when moving towards the decision boundary.

This simple yet effective idea has been proposed for learning methods but has not been applied in the context of digital watermarking. In fact, existing watermarking schemes mainly focus on discriminating marked from unmarked signals and neglect how these are distributed in the media space, leading to the same weak spot. An adversary can therefore exploit the full media space to trigger varying reactions to the respective inputs in order to collect information about the watermark. Broadly speaking, “the image does not have to look nice” in an attack [63]. The so-created points do not necessarily resemble a meaningful signal, but the detector still accepts these points. The Blind Newton Sensitivity Attack from Section 3.3, for instance, needs to find a starting position on the decision boundary. Without further information

about the boundary’s location, an attacker can thus perform a line search in a random direction or set pixels to gray iteratively (see Figure 5 for the resulting images).

Consequently, the defense from the learning community provides us with a new research direction to tackle oracle attacks in digital watermarking. Figure 4(b) depicts a possible *1.5-class watermark detector* that works as follows: The two-class detector enables us to distinguish unmarked from watermarked signals, while the one-class detector enables spotting implausible signals, that is, too far away from reasonable variations. In particular, the two-class detector decides on the presence of a watermark only if the input lies within the one-class region. Outside, the detector returns random decisions or alternatively blocks access for subsequent queries from the same source (see Section 3.5).

To model plausible signals, the 1.5-class watermark detector is trained with samples of honest variations of the target image, such as strong changes of the brightness, compression or denoising. Figure 5 depicts possible variations with distortions where the detector should still decide on watermark presence.

Experimental setup. To demonstrate the practical utility of this novel defense, we conduct an empirical evaluation with a state-of-the-art oracle attack. Our dataset for this evaluation consists of images from the publicly available Dresden Image Database [24], where 50 uncompressed Adobe Lightroom images from a Nikon D70 camera are used. All images are converted to grayscale and cropped to a common size of 128×128 pixels, that is, $N = 16834$ dimensions.

Our experimental procedure is as follows: The watermark embedding and detection process follow the presented scheme in Appendix A that yields a linear decision boundary and that represents the two-class detector. To obtain the respective one-class detector, we create different variations of the watermarked image \tilde{x} by applying common image processing steps, such as noise addition, denoising, JPEG compression as well as contrast-

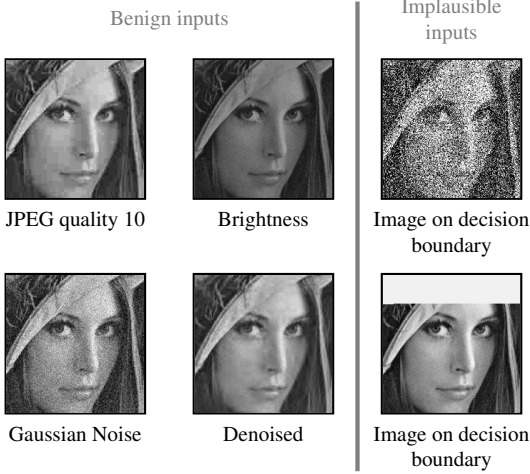


Figure 5: Distortions of the target image. The left four plots show plausible image distortions. The right plots depict boundary starting positions for an attack where the one-class detector gives an alarm.

and brightness variation. We apply neighborhood-based anomaly detection to define a simple model of normality. Given an image \tilde{x} , this model computes the distance d to the k -nearest variation of \tilde{x} , that is,

$$d(\tilde{x}) = \frac{1}{k} \sum_{v \in N_{\tilde{x}}} \|v - \tilde{x}\| \quad (3)$$

where $N_{\tilde{x}}$ are the k -nearest neighbors of \tilde{x} . We mark an image as implausible if the distance to its k -nearest variations reaches a given threshold δ . For our study, we simply fix $k = 3$. We attack our 1.5-class detector using the well-studied Blind Newton Sensitivity Attack [3, 12, 13] that successfully defeats several existing defenses (see Appendix B). We perform the attack against each of the 50 images and report aggregated results.

Benign vs. attack images. We start with a comparison between benign and attack images. For each benign image \tilde{x} , we randomly split its set of image variations into a known partition (75%) and unknown partition (25%). We repeat this procedure 50 times and report the distance of each distorted image from the unknown partition to the one-class model defined by the known partition. Figure 6(a) shows the average distance for each benign image and repetition as well as the attack queries per image \tilde{x} .

The distances between benignly distorted and adversely crafted images are well separable. Without further information about the boundary, an attacker needs to use an arbitrary localized boundary position. Thus, the gradient calculation around that position leads to a large fraction of queries that do not resemble a meaningful signal and exhibit an abnormally high distance. This in turn allows a one-class detector to differentiate between benign and attack inputs.

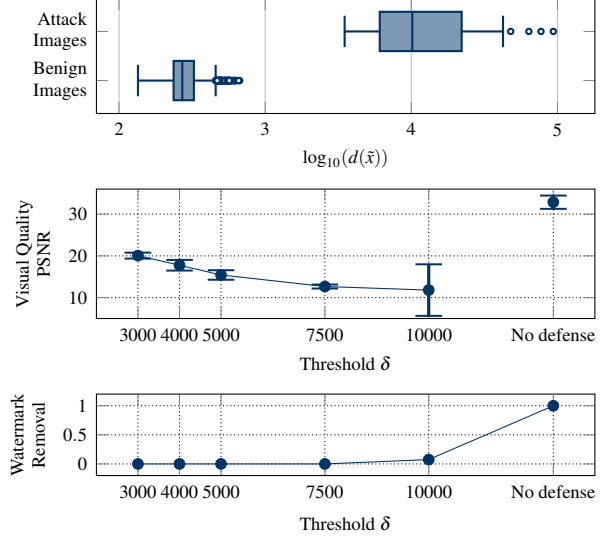


Figure 6: The upper plot compares the average distance from normally distorted and adversely crafted images to a respective one-class model. The middle plot shows the average PSNR of the 50 image outcomes of the attack to the respective original, unwatermarked image as a function of δ . The error bars depict the respective standard deviation. The lower plot depicts the percentage of the attack images where the watermark is not detectable anymore.

Random Decisions. We finally evaluate the impact of random decisions outside the one-class model. Figure 6(b) plots the Peak Signal To Noise Ratio (PSNR) and Figure 6(c) shows the percentage of successfully eliminated watermarks from the 50 final attack outcomes as a function of threshold δ .

A threshold between 3000 and 10000 successfully distorts the attack, such that the watermark is still detectable in the final attack outcome and the visual quality decreases substantially. A higher threshold, however, increases the chances that the attack queries remain inside the one-class model, so that an adversary can render the watermark undetectable with a high visual quality again. A smaller threshold increases the risk of false positives. Moreover, if the model becomes too tight, the attack terminates around the one-class boundary instead. This leads to increasing PSNR values and an adversary could start to exploit the one-class boundary instead. Overall, the results confirm that a suitable threshold strongly impacts the attack with regard to visual quality and watermark removal, and at the same time reduces the risk of false positives.

We acknowledge that false positives may occur with image distortions that are not considered in the one-class model. Yet, the model represents an additional information source for a watermark detector. To the best of our knowledge, a similar defense strategy has not been proposed for digital watermarking so far and in combination with already existing defense mechanisms, we further raise the bar for oracle and watermark-estimation attacks.

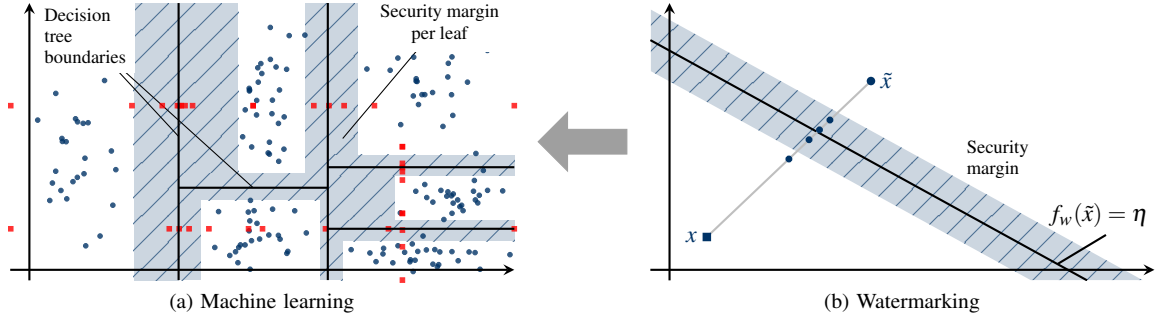


Figure 7: Transfer from digital watermarking to machine learning. The right plot illustrates the concept of the security margin that spots a binary search trying to enclose the boundary. The left plot shows its application to a decision tree. The region defined by each leaf has a security margin where its width is adapted to the training data distribution (circles). In contrast, the queries from the tree extraction algorithm (squares) underline that the attack needs to operate in the security margin to localize the decision boundary.

4.2 From Watermarking to Machine Learning

We proceed with applying concepts from the field of watermarking to machine learning. In particular, we demonstrate that the concept of a stateful detector mitigates the risk of model stealing by identifying sequences of malicious queries. This provides online services offering machine learning as a service new capabilities for fending off attacks. To begin with, we shortly summarize the model extraction proposed by Tramèr et al. [56] for a decision tree and then develop a stateful classifier as an effective countermeasure. Finally, we present an empirical evaluation to demonstrate the practical feasibility of this novel defense strategy.

Decision tree extraction. Tramèr et al. [56] reconstruct decision trees by performing targeted queries on the APIs provided by the BigML service. The attack is possible, since the service does not only return the class label for a submitted query but also a confidence score for a particular leaf node. This enables an adversary to distinguish between the leaves. For each leaf and for each of its features, a recursive binary search locates the leaf’s decision boundary in that direction. As the binary search covers the whole feature range, other leaf regions are discovered as well and extracted subsequently. In this way, an adversary can extract all possible paths of the decision tree. Note that the binary search needs to fix all features except for the one of interest, as otherwise the attack may miss a leaf during the reconstruction.

Stateful decision tree. As a countermeasure to this attack, we devise a defense that builds on a successful protection technique from digital watermarking—a *stateful detector* [1, 55]. Figure 7(b) shows the concept as proposed by Barni et al. [1]. A narrow stripe across the decision boundary determines a *security margin*. The

detector does not only check for the presence of a watermark, but simultaneously counts the number of queries falling inside this margin. An attacker performing a binary search to enclose the boundary will necessarily create an unusually large number of queries in the security margin. The analysis of the input sequences therefore allows the identification of unusual activity which mitigates the risk of oracle and watermark-estimation attacks. The exact parameters of the security margin are derived through statistical properties of the decision function [1]. Although this defense strategy has been initially designed to protect watermark detectors, we demonstrate that it can be extended to secure decision trees as well.

Figure 7(a) illustrates the transferred concept where security margins are added to the boundaries of each tree region. The width of these margins is determined for each region and feature dimension separately depending on the statistical distribution of the data. Overall, the security margin is defined alongside the original decision tree and does not require changes to its implementation. Appendix C provides more information on the margin’s creation process.

When the decision tree returns the predicted class for a query, the stateful detector checks whether the query falls inside the security margin. To determine whether the tree is subject to an attack, we calculate a simple ratio: For each leaf, we count the number of incoming queries. At the same time, the leaf keeps record of the queries inside the security margin. We denote by φ the ratio from the security margin queries to the total number of queries, averaged over all leaves. This ratio is an indicator for the plausibility of the current input sequence. Figure 7(a) also shows the typical query sequence from the tree extraction algorithm (red squared). The adversary has to work in the margin to localize the decision boundary, in contrast to the distribution of benign queries.

Dataset	Original Attack		Blocking Defense		Random Resp. Defense		Adapted Attack	
	Q	p	Q	p	Q	p	Q	p
Iris	108	1.00	38	0.09	*	0.09	4,412	1.00
Carseats	871	1.00	148	0.20	*	0.20	15,156	0.46
College	2,216	1.00	244	0.10	*	0.10	8,974	0.08
Orange Juice	4,804	1.00	846	0.20	*	0.20	86,354	0.48
Wine Quality	9,615	1.00	978	0.11	*	0.11	37,406	0.11

Table 2: Effectiveness of the Security Margin Defense for different attack variations and possible reactions after detecting the attack. Q denotes the number of queries, p the percentage of successfully extracted leaves. Without any defense, the original extraction algorithm from Tramèr et al. extracts the whole tree ($p = 1$). In contrast, with the security margin, the detector can spot an attack and block further access before the whole tree is reconstructed ($p \leq 0.2$). If the adversary adapts the attack by sending cover queries with random values, the attack chances increases, but the full reconstruction is still not possible for larger datasets.

Experimental setup. To evaluate this defense in practice, we use the publicly available tree-stealing implementation by Tramèr et al. [56]. Table 3 summarizes our used datasets. We divide each dataset into a training set (50%) and test set (50%), where we use the first for learning a decision tree and calibrating the security margins. We repeat this process 5 times and present aggregated results in the following. The detector assumes an attack if the query ratio ϕ exceeds the threshold $\tau = 0.3$, estimated by a prior cross-validation.

Dataset	Samples	Features	ϕ Leaves
Iris	150	4	4.6
Carseats	400	8	13.2
College	777	17	18.8
Orange Juice	1,070	11	59.0
Wine Quality	1,599	11	89.4

Table 3: Dataset for evaluation. The number of leaves from the learned decision tree are averaged over the repetitions.

Defense Evaluation. We first examine the security margin under benign and attack queries, where Table 2 reports the results for the corresponding experiments. In the first step, we make use of the test set to simulate the queries of an honest user. In this way, we can determine the risk of false positives, that is, declaring that an honest input sequence is malicious. The final query ratio ϕ after submitting the complete sequence was not higher than 0.2 in all datasets, so that the stateful detector does not mark a benign query sequence as attack by mistake.

Next, we run the tree-stealing attack against the learned tree without and with the security margin defense. In the latter case, we consider two reactions after that an attack sequence is detected: (a) the tree blocks further access, (b) the tree returns random decisions. To determine the knowledge gain by the adversary, Table 2 reports the percentage of successfully extracted leaves p . The blocking strategy allows the tree to block the tree extraction at the very beginning. With random decisions, the attack’s binary search recursively locates an exponential number

of boundaries erroneously. We stopped the attack after 1 Million queries (marked by *).

Counter-Attack Evaluation. As a counter-reaction, an adversary can in turn submit *cover queries* outside the security margin so that the query ratio ϕ ideally remains below the threshold. There are, however, two practical problems. Without knowledge of the training data distribution, the adversary cannot know where a decision boundary could be located and thus where the margin could be. Another problem is that the attacker needs to control the ratio in almost each leaf. It is not sufficient to send just one fixed well-chosen cover query all the time, since this query would only affect one leaf. These two problems complicate the design of cover queries.

We therefore let the attacker create cover queries by selecting random values in the range of each feature. Table 2 shows the performance of this adapted attack where an adversary sends 40 cover queries for each tree extraction query. Still, the whole tree cannot be extracted. Only half of the leaves are extracted before the detector spots the attack and blocks further access.

We finally consider a stronger attacker who knows a certain percentage of the training data. This is not unrealistic, if an adversary can make some assumptions about possible training data. The attacker can make use of the leaked training data as cover queries. Table 4 summarizes the percentage of extracted leaves p for varying amounts of known training data and cover queries. If just 10% of the data are known, even 40 cover queries between each attack query do not suffice to extract the whole tree. However, if the adversary knows more data points, the cover queries spread over all leaves more equally and the attack chances start to increase.

Overall, our evaluation demonstrates that the proposed defense can effectively mitigate the risk of model stealing based on the history of queries. While this defense is only a first step in hindering model-extraction attacks, we show that the concept of a stateful analysis brings in a new defense strategy that in combination with other protections,

Dataset	Cover Queries	Percentage train. data				
		10	20	30	40	50
Iris	1x	0.17	0.21	0.21	0.21	0.22
	5x	0.64	0.85	0.89	0.92	0.94
	40x	0.76	0.91	0.94	0.97	1.00
Carseats	1x	0.28	0.29	0.28	0.29	0.30
	5x	0.39	0.60	0.69	0.82	0.89
	40x	0.50	0.87	0.97	1.00	1.00
College	1x	0.12	0.12	0.12	0.12	0.12
	5x	0.17	0.26	0.28	0.29	0.32
	40x	0.29	0.64	0.85	0.94	1.00
Orange Juice	1x	0.28	0.29	0.29	0.29	0.29
	5x	0.39	0.63	0.88	0.98	0.99
	40x	0.46	0.92	1.00	1.00	1.00
Wine Quality	1x	0.20	0.22	0.22	0.23	0.24
	5x	0.33	0.55	0.88	0.98	1.00
	40x	0.43	0.91	1.00	1.00	1.00

Table 4: Percentage of extracted leaves with an informed attacker who knows a certain percentage of the training data. Results are shown for different numbers of cover queries that are sent between each attack query from the tree extraction algorithm.

such as a line search detection as proposed for watermarking [1], can lower the chances of reconstructing a model in reasonable time. As our defense can be implemented alongside an existing classifier, online services such as BigML can easily deploy our defense in practice.

5 Discussion

Adversarial machine learning and digital watermarking are vivid research fields that have established a broad range of methods and concepts. However, one can observe the following asymmetry: the former community has focused on white-box attacks, while the latter has extensively studied the black-box threat. Although recent research in machine learning has started to also study black-box attacks more thoroughly [39, 40, 56], existing insights from digital watermarking potentially bring forth novel ideas. The other way round, we also show that knowledge from machine learning can help to mitigate attacks against watermarks.

The presented unified view opens interesting directions for future work. The comparison of defenses between both fields in Section 3.5, for instance, discloses that stateful detection strategies have not been considered in machine learning yet. While we successfully transferred the security margin approach in this paper, a line search detection based on a PCA, for example, could further mitigate model-extraction attacks. On the other side, the various strategies, such as adaptive re-learning [56], that have been used successfully against classifiers mark a potential threat for watermark detectors as well.

Furthermore, the watermarking community has conducted different contests, where researchers could attack and defend watermarking schemes under practical

conditions, such as the “Break Our Watermarking System” (BOWS) competition. These contests have promoted a variety of publications that reveal shortcomings of existing protection techniques and introduce novel defenses [e.g. 4, 13, 16, 61, 62, 64], such as the strong watermarking scheme *Broken Arrows* [21]. Based on our unified view, we encourage the organization of a similar *contest for adversarial machine learning*. By imposing researchers into the role of an attacker in a real scenario without perfect knowledge, previously unknown questions and insights often come to light. Šrndić and Laskov [59], for instance, demonstrate the feasibility to evade a publicly available PDF malware classifier with the insight that full knowledge of the classifier features is not necessary. The contest could be structured similarly to watermarking contests in different episodes, with each providing a different level of knowledge about a defense or attack (see Appendix D).

Finally, we note that machine learning and watermarking are not the only research areas that have to cope with an adversary. The identified similarities between both research fields can be seen as part of a bigger problem: *Adversarial Signal Processing* [2]. More fields such as information or multimedia forensics also deal with an adversary’s presence and to our knowledge the transfer of concepts between and to these areas has not been addressed so far.

6 Conclusion

Developing analysis methods for an adversarial environment is a challenging task: First, these methods need to provide correct results even if parts of their input are manipulated and, second, these methods should protect from known as well as future attacks. The research fields of adversarial learning and digital watermarking both have tackled this challenge and developed a remarkable set of defenses for operating in an adversarial environment.

In this paper, we show that both lines of research share similarities which have been overlooked by previous work and enable transferring concepts from one field to the other. By means of a systematization of attacks, we are able to transform defenses for learning methods to the domain of watermarking and vice versa. This not only opens new perspective for designing joint defenses, but also allows for combining techniques from both fields that have not been previously coupled.

As part of our analysis, we identify interesting directions of future research that enable the two communities to learn from each other and combine the “best of both worlds”. As one example of these directions, we particularly encourage the organization of a public competition for adversarial machine learning, where attacks and defenses are put to the test in a competitive manner.

References

- [1] BARNI, M., COMESAÑA-ALFARO, P., PÉREZ-GONZÁLEZ, F., AND TONDI, B. Are you threatening me?: Towards smart detectors in watermarking. *Proceedings of SPIE 9028* (2014).
- [2] BARNI, M., AND PÉREZ-GONZÁLEZ, F. Coping with the enemy: Advances in adversary-aware signal processing. In *IEEE International Conference on Acoustics, Speech, and Signal Processing*. 2013, pp. 8682–8686.
- [3] BARNI, M., PÉREZ-GONZÁLEZ, F., COMESAÑA, P., AND BARTOLI, G. Putting reproducible signal processing into practice: A case study in watermarking. In *International Conference on Acoustics, Speech and Signal Processing (ICASSP)* (2007), pp. 1261–1264.
- [4] BAS, P., AND WESTFELD, A. Two key estimation techniques for the broken arrows watermarking scheme. In *Proc. of ACM Workshop on Multimedia and Security* (2009), pp. 1–8.
- [5] BBC. Pornographic films on BitTorrent: Flava Works gets huge damages. <http://www.bbc.co.uk/news/technology-20178171>. last visited October 2016.
- [6] BIGGIO, B., CORONA, I., HE, Z., CHAN, P. P. K., GIACINTO, G., YEUNG, D. S., AND ROLI, F. One-and-a-half-class multiple classifier systems for secure learning against evasion attacks at test time. In *Proc. of International Workshop on Multiple Classifier Systems (MCS)* (2015).
- [7] BIGGIO, B., CORONA, I., MAIORCA, D., NELSON, B., ŠRNDIĆ, N., LASKOV, P., GIACINTO, G., AND ROLI, F. Evasion attacks against machine learning at test time. In *Machine Learning and Knowledge Discovery in Databases*. Springer, 2013, pp. 387–402.
- [8] BIGGIO, B., FUMERA, G., AND ROLI, F. Adversarial pattern classification using multiple classifiers and randomisation. In *Structural, Syntactic, and Statistical Pattern Recognition*. Springer, 2008, pp. 500–509.
- [9] BIGGIO, B., NELSON, B., AND LASKOV, P. Support vector machines under adversarial label noise. In *Proc. of Asian Conference on Machine Learning (ACML)* (2011), pp. 97–112.
- [10] BIGGIO, B., NELSON, B., AND LASKOV, P. Poisoning attacks against support vector machines. In *Proc. of International Conference on Machine Learning (ICML)* (2012).
- [11] CHOUBASSI, M. E., AND MOULIN, P. Noniterative algorithms for sensitivity analysis attacks. *IEEE Transactions on Information Forensics and Security* 2, 2 (2007), 113–126.
- [12] COMESAÑA, P., PÉREZ-FREIRE, L., AND PÉREZ-GONZÁLEZ, F. Blind newton sensitivity attack. *IEE Proceedings – Information Security* 153, 3 (2006), 115–125.
- [13] COMESAÑA, P., AND PÉREZ-GONZÁLEZ, F. Breaking the BOWS watermarking system: Key guessing and sensitivity attacks. *EURASIP Journal on Information Security* 2007, 1 (2007).
- [14] COX, I. J., AND LINNARTZ, J.-P. M. G. Public watermarks and resistance to tampering. In *Proc. of IEEE International Conference on Image Processing (ICIP)* (1997), pp. 26–29.
- [15] COX, I. J., MILLER, M., BLOOM, J., FRIDRICH, J., AND KALKER, T. *Digital watermarking and steganography*. Morgan Kaufmann Publishers, 2002.
- [16] CRAVER, S., AND YU, J. Reverse-engineering a detector with false alarms. *Proceedings of SPIE 6505* (2007), 65050C.
- [17] DUDA, R., P.E.HART, AND D.G.STORK. *Pattern classification*, second ed. John Wiley & Sons, 2001.
- [18] EL CHOUBASSI, M., AND MOULIN, P. On the fundamental tradeoff between watermark detection performance and robustness against sensitivity analysis attacks. *Proceedings of SPIE 6072* (2006), 1–12.
- [19] FOGLA, P., AND LEE, W. Evading network anomaly detection systems: formal reasoning and practical techniques. In *Proc. of ACM Conference on Computer and Communications Security (CCS)* (2006), pp. 59–68.
- [20] FOGLA, P., SHARIF, M., PERDISCI, R., KOLESNIKOV, O., AND LEE, W. Polymorphic blending attacks. In *Proc. of USENIX Security Symposium* (2006), pp. 241–256.
- [21] FURON, T., AND BAS, P. Broken arrows. *EURASIP Journal on Information Security* 2008 (2008), 1–13.
- [22] FURON, T., MACQ, B., HURLEY, N., AND SILVESTRE, G. JANIS: Just another N-order side-informed watermarking scheme. In *Proc. of IEEE International Conference on Image Processing (ICIP)* (2002), vol. 3, pp. 153–156.
- [23] FURON, T., VENTURINI, I., AND DUHAMEL, P. Unified approach of asymmetric watermarking schemes. *Proceedings of SPIE 4314* (2001), 269–279.
- [24] GLOE, T., AND BÖHME, R. The Dresden Image Database for benchmarking digital image forensics. *Journal of Digital Forensic Practice* 3, 2–4 (2010), 150–159.
- [25] GROSSE, K., PAPERNOT, N., MANOHARAN, P., BACKES, M., AND MCDANIEL, P. Adversarial perturbations against deep neural networks for malware classification. Tech. Rep. abs/1606.04435, Computing Research Repository (CoRR), 2016.
- [26] HASTIE, T., TIBSHIRANI, R., AND FRIEDMAN, J. *The Elements of Statistical Learning: data mining, inference and prediction*. Springer series in statistics. Springer, New York, N.Y., 2001.
- [27] HUANG, L., JOSEPH, A. D., NELSON, B., RUBINSTEIN, B. I. P., AND TYGAR, J. D. Adversarial machine learning. In *Proc. of ACM Workshop on Artificial Intelligence and Security (AISEC)* (2011), pp. 43–58.
- [28] KALKER, T., LINNARTZ, J.-P. M. G., AND VAN DIJK, M. Watermark estimation through detector analysis. In *Proc. of IEEE International Conference on Image Processing (ICIP)* (1998), pp. 425–429.
- [29] KAPRAVELOS, A., SHOSHITAISHVILI, Y., COVA, M., KRUEGEL, C., AND VIGNA, G. Revolver: An automated approach to the detection of evasive web-based malware. In *Proc. of USENIX Security Symposium* (Aug. 2013), pp. 637–651.
- [30] KOLBITSCH, C., LIVSHITS, B., ZORN, B., AND SEIFERT, C. Rozzle: De-cloaking internet malware. In *Proc. of IEEE Symposium on Security and Privacy* (2012), pp. 443–457.
- [31] KOŁCZ, A., AND TEO, C. H. Feature weighting for improved classifier robustness. In *Proc. of Conference on Email and Anti-Spam (CEAS)* (2009).
- [32] LEVINSON, J., ASKELAND, J., BECKER, J., DOLSON, J., HELD, D., KAMMEL, S., KOLTER, J. Z., LANGER, D., PINK, O., PRATT, V., SOKOLSKY, M., STANEK, G., STAVENS, D. M., TEICHMAN, A., WERLING, M., AND THRUN, S. Towards fully autonomous driving: Systems and algorithms. In *Proc. of IEEE Intelligent Vehicles Symposium (IV)* (2011), pp. 163–168.

- [33] LIAO, X., YUAN, K., WANG, X., LI, Z., XING, L., AND BEYAH, R. A. Acing the ioc game: Toward automatic discovery and analysis of open-source cyber threat intelligence. In *Proc. of ACM Conference on Computer and Communications Security (CCS)* (2016), pp. 755–766.
- [34] LINNARTZ, J.-P. M. G., AND VAN DIJK, M. Analysis of the sensitivity attack against electronic watermarks in images. In *Proc. of Information Hiding Conference* (1998), vol. 1525, pp. 258–272.
- [35] LOWD, D., AND MEEK, C. Adversarial learning. In *Proc. of ACM SIGKDD Conference on Knowledge Discovery in Data Mining (KDD)* (2005), pp. 641–647.
- [36] LOWD, D., AND MEEK, C. Good word attacks on statistical spam filters. In *Conference on Email and Anti-Spam* (2005).
- [37] MANSOUR, M. F., AND TEWFIK, A. H. Improving the security of watermark public detectors. In *Proc. of International Conference on Digital Signal Processing (DSP)* (2002), pp. 59–66.
- [38] MANSOUR, M. F., AND TEWFIK, A. H. LMS-based attack on watermark public detectors. In *Proc. of International Conference on Image Processing (ICIP)* (2002), pp. 649–652.
- [39] PAPERNOT, N., MCDANIEL, P., AND GOODFELLOW, I. Transferability in machine learning: from phenomena to black-box attacks using adversarial samples. Tech. rep., arXiv:1605.07277, 2016.
- [40] PAPERNOT, N., MCDANIEL, P., GOODFELLOW, I., JHA, S., BERKAY CELIK, Z., AND SWAMI, A. Practical black-box attacks against deep learning systems using adversarial examples. Tech. rep., arXiv:1602.02697, 2016.
- [41] PAPERNOT, N., MCDANIEL, P., JHA, S., FREDRIKSON, M., CELIK, Z. B., AND SWAMI, A. The limitations of deep learning in adversarial settings. In *Proc. of IEEE European Symposium on Security and Privacy* (2016).
- [42] PAPERNOT, N., MCDANIEL, P., SINHA, A., AND WELLMAN, M. Towards the science of security and privacy in machine learning. Tech. Rep. abs/1611.03814, Computing Research Repository (CoRR), 2016.
- [43] PAPERNOT, N., MCDANIEL, P. D., WU, X., JHA, S., AND SWAMI, A. Distillation as a defense to adversarial perturbations against deep neural networks. In *Proc. of IEEE Symposium on Security and Privacy* (2016).
- [44] PERDISCI, R., GU, G., AND LEE, W. Using an ensemble of one-class SVM classifiers to harden payload-based anomaly detection systems. In *Proc. of International Conference on Data Mining (ICDM)* (2006), pp. 488–498.
- [45] PIVA, A., AND BARNI, M. Design and analysis of the first BOWS contest. *EURASIP Journal on Information Security 2007* (2007), 3:1–3:7.
- [46] PROVOS, N., AND HONEYMAN, P. Detecting steganographic content on the internet. In *Proc. of Network and Distributed System Security Symposium (NDSS)* (2002).
- [47] RUSSU, P., DEMONTIS, A., BIGGIO, B., FUMERA, G., AND ROLI, F. Secure kernel machines against evasion attacks. In *Proc. of ACM Workshop on Artificial Intelligence and Security (AISEC)* (2016), pp. 59–69.
- [48] SCHÖLKOPF, B., AND SMOLA, A. J. *Learning with Kernels*. MIT Press, Cambridge, MA, 2002.
- [49] SCHROFF, F., KALENICHENKO, D., AND PHILBIN, J. Facenet: A unified embedding for face recognition and clustering. In *Proc. of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2015), pp. 815–823.
- [50] SHARIF, M., BHAGAVATULA, S., BAUER, L., AND REITER, M. K. Accessorize to a crime: Real and stealthy attacks on state-of-the-art face recognition. In *Proc. of ACM Conference on Computer and Communications Security (CCS)* (2016), pp. 1528–1540.
- [51] SHOKRI, R., STRONATI, M., AND SHMATIKOV, V. Membership inference attacks against machine learning models. Tech. rep., arXiv:1610.05820v1, 2016.
- [52] SONG, Y., LOCASO, M., STAVROU, A., KEROMYTIS, A., AND STOLFO, S. On the infeasibility of modeling polymorphic shellcode. In *Proc. of ACM Conference on Computer and Communications Security (CCS)* (2007), pp. 541–551.
- [53] SRNDIC, N., AND LASKOV, P. Practical evasion of a learning-based classifier: A case study. In *Proc. of IEEE Symposium on Security and Privacy* (2014), pp. 197–211.
- [54] TAIGMAN, Y., YANG, M., RANZATO, M. A., AND WOLF, L. Deepface: Closing the gap to human-level performance in face verification. In *Proc. of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2014).
- [55] TONDI, B., COMESAÑA-ALFARO, P., PÉREZ-GONZÁLEZ, F., AND BARNI, M. On the effectiveness of meta-detection for countering oracle attacks in watermarking. In *Workshop on Information Forensics and Security (WIFS)* (2015), pp. 1–6.
- [56] TRAMÈR, F., ZHANG, F., JUELS, A., REITER, M. K., AND RISTENPART, T. Stealing machine learning models via prediction apis. In *Proc. of USENIX Security Symposium* (2016), pp. 601–618.
- [57] VENKATESAN, R., AND JAKUBOWSKI, M. H. Randomized detection for spread-spectrum watermarking: Defending against sensitivity and other attacks. In *Proc. of IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)* (2005), vol. 2, pp. 9–12.
- [58] VENTURINI, I. Oracle attacks and covert channels. In *Proc. of International Workshop on Digital Watermarking* (2005), vol. 3710, pp. 171–185.
- [59] ŠRNDIĆ, N., AND LASKOV, P. Practical evasion of a learning-based classifier: A case study. In *Proc. of IEEE Symposium on Security and Privacy* (2014).
- [60] WANG, K., PAREKH, J. J., AND STOLFO, S. J. Anagram: A content anomaly detector resistant to mimicry attack. In *Proc. of International Symposium on Recent Advances in Intrusion Detection (RAID)* (2006), pp. 226–248.
- [61] WESTFELD, A. Lessons from the BOWS contest. In *Workshop on Multimedia and Security (MM&Sec)* (2006), pp. 208–213.
- [62] WESTFELD, A. A workbench for the BOWS contest. *EURASIP Journal on Information Security 2007*, 1 (2008), 064521.
- [63] WESTFELD, A. Fast determination of sensitivity in the presence of countermeasures in BOWS-2. In *International Workshop on Information Hiding*. Springer, 2009, pp. 89–101.
- [64] XIE, F., FURON, T., AND FONTAINE, C. Better security levels for broken arrows. *Proceedings of SPIE 7541* (2010), 75410H.

- [65] XU, W., QI, Y., AND EVANS, D. Automatically evading classifiers: A case study on pdf malware classifiers. In *Proc. of Network and Distributed System Security Symposium (NDSS)* (2016).
- [66] ZHU, Z., LIANG, D., ZHANG, S., HUANG, X., LI, B., AND HU, S. Traffic-sign detection and classification in the wild. In *Proc. of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2016).

A Linear Watermark Detector

This section illustrates the watermarking process with a simple watermarking scheme. The process can be generally divided into two phases: embedding and detection. Let us, for instance, consider the *additive spectrum watermarking scheme* which is also the embedding scheme used in the image example from Figure 2. In this scheme, the watermarked version \tilde{x} of a signal x is created by adding a watermarking vector $w \in \mathbb{R}^N$ onto x element-wise, that is,

$$\tilde{x} = x + w. \quad (4)$$

The watermark w usually represents a random pattern. In order to decide whether a signal contains the particular watermark, a *linear correlation detector* can be employed that uses the following decision function

$$f_w(\tilde{x}) = \tilde{x}^\top w. \quad (5)$$

The output is a weighted sum between \tilde{x} and the watermark w . If watermark and signal match, the correlation exceeds a pre-defined threshold η . Geometrically, each signal corresponds to a point in a vector space where the watermark describes a decision boundary, as shown in Figure 8. The result are two subspaces, one for the watermark's presence, one for its absence. The detection thus works by determining which subspace an input signal is currently in.

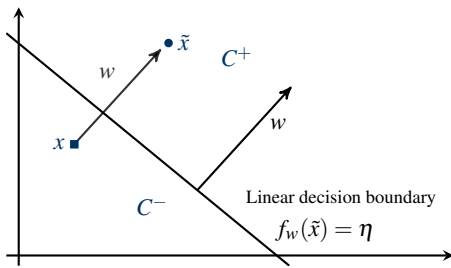


Figure 8: Geometrical view on the embedding and detection process of a simple watermarking scheme.

B Blind Newton Sensitivity Attack

This section briefly recaps the Blind Newton Sensitivity Attack (BNSA) that interprets the watermark removal as

a non-linear optimization problem [12]:

$$\min d(t) \quad (6)$$

$$\text{subject to } f_w(\tilde{x}') = f_w(\tilde{x} + t) = \eta. \quad (7)$$

The objective function $d(t)$ measures the changes t on the image \tilde{x} . For example, the squared Euclidean norm $d(t) = \|t\|_2^2$ minimizes the length of t and therefore the necessary pixel changes. At the same time, the optimal solution must satisfy the constraint that the detector does not detect the watermark. A position on the boundary is here sufficient, so that Equation (7) restricts the decision function to η .

The adversary, however, does not know the function $f_w(\tilde{x}')$, since the watermark w is kept secret. Although only a binary detector output is observable, an attack is yet possible. To this end, Comesaña et al. rewrite the optimization problem into an unconstrained version:

$$\arg \min_{t \in \mathbb{R}^N} d(h(t)). \quad (8)$$

The function $h(t)$ reflects the prior constraint by mapping t to the decision boundary. To this end, a bisection algorithm can be used to find a scalar α such that αt lies on the decision boundary. There is, however, another problem. As $h(t)$ has to map each input vector to the boundary explicitly by running the bisection algorithm respectively, a closed form to solve the problem is not applicable. Therefore numeric iterative methods such as Newton's method or gradient descent have to be used as Figure 9 exemplifies.

The attack starts with a random direction to locate the decision boundary. After calculating an image at the boundary, it slightly changes the vector at one position, maps the vector to the boundary again and records the distance through this change. By repeating this procedure for each feature direction, the attack is able to calculate the gradient at this boundary position. This step yields the direction in which the necessary changes $d(t)$ decreases fastest. In this way, the attack is able to locate a boundary position that is closer to \tilde{x} than the previous position. This process can be repeated, but in the case of a linear boundary, the algorithm finishes after one iteration with the fewest necessary changes.

In summary, the attack does not require a priori knowledge about the detector's decision function and works only with a binary output. Although the attack only converges to an optimal solution for convex boundaries, it has been used against various watermarking schemes with even polynomial and fractalized decision boundaries [12, 13, 55].

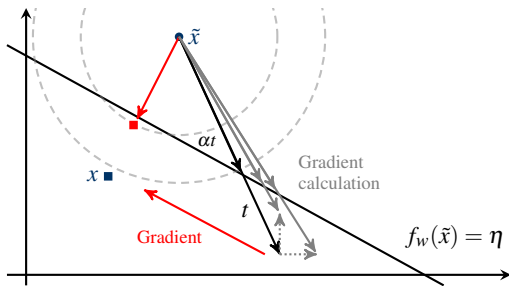


Figure 9: Blind Newton Sensitivity Attack. Queries around a boundary position reveal the function’s gradient at this position to minimize the distance between the manipulated sample and the original one.

C Security Margin Construction

The security margin’s construction works as follows: First, we choose a tree region and select the training data that fall inside this particular region. Next, we estimate the distribution of the selected training data at each dimension through a kernel-density estimation. In this way, no a priori assumptions about their distribution are required. Finally, the distribution in each dimension is used to define the margin at the boundary in this dimension. To this end, we set the margin to the feature value where the probability of occurrence is smaller than a certain threshold. In Figure 7(a), for example, the top right tree regions has a smaller security margin, since more training data are near the boundary. On the contrary, the most left region exhibits fewer training samples near the boundary, so that a larger margin can be defined. By defining the security margin in this statistical way, we can control the false alarm rate that a honest query falls inside the margin. We repeat the process for each tree region.

D BOWS Contest

“Break Our Watermarking System” or BOWS is a contest that has been held twice in the watermarking community. The latest contest is divided into three subsequent episodes, where only the last episode reveals the underlying watermarking scheme.

1. At the beginning of the contest, 3 watermarked images are available together with an online watermark detector that allows 30 calls per day. This episode models an attacker with limited knowledge and capabilities. The participants are required to operate with few queries and need to carefully construct their attacks.
2. In the next episode, the daily rate limited is dropped and the participants can perform different forms of oracle and watermark-estimation attacks against the detector. The episode models a stronger attacker, yet

only 3 images are available for inferring the pattern embedded by the watermarking scheme.

3. Finally, the same watermark is embedded into 10,000 images and the underlying watermarking scheme is released. This episode models a very strong adversary with full knowledge of the scheme together with access to a large set of images. Ultimately, a watermarking scheme should remain secure even in this setting.

For each episode, a hall of fame on the respective website documents the participant’s success regarding the image quality and the launched attacks. Further information on the design of both contests are provided by Piva and Barni [45] as well as Furon and Bas [21], respectively.