

Randomized Prediction Games for Adversarial Machine Learning

Samuel Rota Bulò, *Member, IEEE*, Battista Biggio, *Member, IEEE*, Ignazio Pillai, Marcello Pelillo, *Fellow, IEEE*, and Fabio Roli, *Fellow, IEEE*

Abstract—In spam and malware detection, attackers exploit randomization to obfuscate malicious data and increase their chances of evading detection at test time, e.g., malware code is typically obfuscated using random strings or byte sequences to hide known exploits. Interestingly, randomization has also been proposed to improve security of learning algorithms against evasion attacks, as it results in hiding information about the classifier to the attacker. Recent work has proposed game-theoretical formulations to learn secure classifiers, by simulating different evasion attacks and modifying the classification function accordingly. However, both the classification function and the simulated data manipulations have been modeled in a deterministic manner, without accounting for any form of randomization. In this paper, we overcome this limitation by proposing a randomized prediction game, namely, a noncooperative game-theoretic formulation in which the classifier and the attacker make randomized strategy selections according to some probability distribution defined over the respective strategy set. We show that our approach allows one to improve the tradeoff between attack detection and false alarms with respect to the state-of-the-art secure classifiers, even against attacks that are different from those hypothesized during design, on application examples including handwritten digit recognition, spam, and malware detection.

Index Terms—Adversarial learning, computer security, evasion attacks, game theory, pattern classification, randomization.

I. INTRODUCTION

MACHINE-LEARNING algorithms have been increasingly adopted in adversarial settings, such as spam, malware, and intrusion detection. However, such algorithms are not designed to operate against intelligent and adaptive attackers, thus making them inherently vulnerable to carefully crafted attacks. Evaluating security of machine learning against such attacks and devising suitable countermeasures are two among the main open issues under investigation in the field of adversarial machine learning [1]–[11]. In this paper, we focus on the issue of designing secure classification algorithms against evasion attacks, i.e., attacks in which malicious samples are manipulated at test time to evade detection. This is a

typical setting, e.g., in spam filtering, where spammers manipulate the content of spam emails to get them past the antispyam filters [1], [2], [12]–[14], or in malware detection, where hackers obfuscate malicious software (malware, for short) to evade detection of either known or zero-day exploits [8], [9], [15], [16]. Although out of the scope of this paper, it is worth mentioning here another pertinent attack scenario, referred to as classifier poisoning. Under this setting, the attacker can manipulate the training data to mislead classifier learning and cause a denial of service, e.g., by increasing the number of misclassified samples [6], [7], [17]–[20].

To date, several authors have addressed the problem of designing secure learning algorithms to mitigate the impact of evasion attacks [1], [6], [10], [11], [21]–[27] (see Section VII for further details). The underlying rationale of such approaches is to learn a classification function that accounts for potential malicious data manipulations at test time. To this end, the interactions between the classifier and the attacker are modeled as a game in which the attacker manipulates data to evade detection, while the classification function is modified to classify them correctly. This essentially amounts to incorporating knowledge of the attack strategy into the learning algorithm. However, both the classification function and the simulated data manipulations have been modeled in a deterministic manner, without accounting for any form of randomization.

Randomization is often used by attackers to increase their chances of evading detection, e.g., malware code is typically obfuscated using random strings or byte sequences to hide known exploits, and spam often contains bogus text randomly taken from English dictionaries to reduce the spamminess of the overall message. Surprisingly, randomization has also been proposed to improve classifier security against evasion attacks [3], [6], [28]. It has been shown that randomizing the learning algorithm may effectively hide information about the classification function to the attacker, requiring her to select a less effective attack (manipulation) strategy. In practice, the fact that the adversary may not know the classification function exactly (i.e., in a deterministic sense) decreases her (expected) payoff on each attack sample; i.e., to achieve the same expected evasion rate attained in the deterministic case, the attacker has to increase the number of modifications made to the attack samples [28].

Motivated by the aforementioned facts, in this paper, we generalize static prediction games, i.e., the game-theoretical formulation proposed in [10] and [11], to account for randomized classifiers and data manipulation strategies. For this

Manuscript received April 9, 2016; accepted July 14, 2016.

S. Rota Bulò is with ICT-Tev, Fondazione Bruno Kessler, Trento 38123, Italy (e-mail: rotabulo@fbk.eu).

B. Biggio, I. Pillai, and F. Roli are with the Department of Electrical and Electronic Engineering, University of Cagliari, Cagliari 09123, Italy (e-mail: battista.biggio@diee.unica.it; pillai@diee.unica.it; roli@diee.unica.it).

M. Pelillo is with the Dipartimento di Scienze Ambientali, Informatica e Statistica, Ca' Foscari University of Venice, Venice 30123, Italy (e-mail: pelillo@dsi.unive.it).

This paper has supplementary downloadable material available at <http://ieeexplore.ieee.org>, provided by the author.

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TNNLS.2016.2593488

reason, we refer to our game as a randomized prediction game. A randomized prediction game is a noncooperative game between a randomized learner and a randomized attacker (also called data generator), where the player's strategies are replaced with probability distributions defined over the respective strategy sets. Our goal is twofold. We do not only aim to assess whether randomization helps achieving a better tradeoff in terms of false alarms and attack detection (with respect to the state-of-the-art secure classifiers), but also whether our approach remains more secure against attacks that are different from those hypothesized during design. In fact, given that our game considers randomized players, it is reasonable to expect that it may be more robust to potential deviations from its original hypotheses about the players' strategies.

This paper is structured as follows. Randomized prediction games are presented in Section II, where sufficient conditions for the existence and uniqueness of a Nash equilibrium in these games are also given. In Section III, we focus on a specific game instance involving a linear support vector machine (SVM) learner, for which we provide an effective method to find an equilibrium by overcoming some computational problems. We discuss how to enable the use of nonlinear (kernelized) SVMs in Section IV. In Section V, we report a simple example to explain how the proposed methods enforce security in adversarial settings. Related work is discussed in Section VII. In Section VI, we empirically validate the soundness of the proposed approach on an handwritten digit recognition task, and on adversarial application examples involving spam filtering and malware detection in PDF files. Notably, to evaluate robustness of our approach and the state-of-the-art secure classification algorithms, we also consider attacks that deviate from the models hypothesized during classifier design. Finally, in Section VIII, we summarize our contributions and sketch potential directions for future work.

II. RANDOMIZED PREDICTION GAMES

Consider an adversarial learning setting involving two actors: a data generator and a learner.¹ The data generator produces at training time a set $\hat{\mathcal{D}} = \{\hat{\mathbf{x}}_i, y_i\}_{i=1}^n \subseteq \mathcal{X} \times \mathcal{Y}$ of n training samples, sampled from an unknown probability distribution. Sets \mathcal{X} and \mathcal{Y} denote the input and output spaces of the learning task, respectively. At test time, the data generator modifies the samples in $\hat{\mathcal{D}}$ to form a new data set $\mathcal{D} \subseteq \mathcal{X} \times \mathcal{Y}$, reflecting a test distribution, which differs in general from the training distribution and it is not available at training time. We assume binary learners, i.e., $\mathcal{Y} = \{-1, +1\}$, and we also assume that the data transformation process leaves the labels of the samples in $\hat{\mathcal{D}}$ unchanged, i.e., $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$. Hence, a perturbed data set will simply be represented in terms of a tuple $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n) \in \mathcal{X}^n$, each element being the perturbation of the original input sample $\hat{\mathbf{x}}_i$, while we implicitly assume the label to remain y_i . The role of the learner is to classify samples $\mathbf{x} \in \mathcal{X}$ according to the prediction function $h(\mathbf{x}) = \text{sign}[f(\mathbf{x}; \mathbf{w})]$, which is expressed in terms of a linear generalized decision function $f(\mathbf{x}; \mathbf{w}) = \mathbf{w}^\top \phi(\mathbf{x})$, where $\mathbf{w} \in \mathbb{R}^m$, $\mathbf{x} \in \mathcal{X}$, and $\phi: \mathcal{X} \rightarrow \mathbb{R}^m$ is a feature map.

Static prediction games have been introduced in [11] by modeling the learner and the data generator as the players of a noncooperative game that we identify as l -player and d -player, respectively. The strategies of l -player correspond to the parameterizations \mathbf{w} of the prediction function f . The strategies of the data generator, instead, are assumed to live directly in the feature space, by regarding $\dot{\mathbf{X}} = (\dot{\mathbf{x}}_1^\top, \dots, \dot{\mathbf{x}}_n^\top)^\top \in \mathbb{R}^{mn}$ as a data generator strategy, where $\dot{\mathbf{x}}_i = \phi(\mathbf{x}_i)$. By doing so, the decision function f becomes linear in either players' strategies. Each player is characterized also by a cost function that depends on the strategies played by either players. The cost function of d -player and l -player is denoted by c_d and c_l , respectively, and is given by

$$c_l(\mathbf{w}, \dot{\mathbf{X}}) = \rho_l \Omega_l(\mathbf{w}) + \sum_{i=1}^n \ell_l(\mathbf{w}^\top \dot{\mathbf{x}}_i, y_i) \quad (1)$$

$$c_d(\mathbf{w}, \dot{\mathbf{X}}) = \rho_d \Omega_d(\dot{\mathbf{X}}) + \sum_{i=1}^n \ell_d(\mathbf{w}^\top \dot{\mathbf{x}}_i, y_i) \quad (2)$$

where $\mathbf{w} \in \mathbb{R}^m$ is the strategy of l -player, $\dot{\mathbf{X}} \in \mathbb{R}^{mn}$ is the strategy of the d -player, and y_i denotes the label of $\dot{\mathbf{x}}_i$ as per $\hat{\mathcal{D}}$. Moreover, $\rho_{d/l} > 0$ is a tradeoff parameter, $\ell_{d/l}(\mathbf{w}^\top \dot{\mathbf{x}}_i, y)$ measures the loss incurred by the l/d -player when the decision function yields $\mathbf{w}^\top \dot{\mathbf{x}}_i$ for the i th training sample while the true label is y , and $\Omega_{d/l}$ can be regarded as a penalization for playing a specific strategy. For the d -player, this term quantifies the cost of perturbing $\hat{\mathcal{D}}$ in feature space.

The goal of this paper is to introduce a randomization component in the model of [11], particularly to what concerns the players' behavior. To this end, we take one abstraction step with respect to the aforementioned prediction game, where we let the learner and the data generator sample their playing strategy in the prediction game from a parameterized distribution, under the assumption that they are expected cost-minimizing (also known as expected utility maximizing). By doing so, we introduce a new noncooperative game that we call randomized prediction game between the l -player and the d -player with strategies being mapped to the possible parameterizations of the players' respective distributions, and cost functions being expected costs under the same distributions.

A. Definition of Randomized Prediction Game

Consider a prediction game as described earlier. We inject randomness in the game by assigning each player a parameterized probability distribution, i.e., $p_l(\mathbf{w}; \boldsymbol{\theta}_l)$ for the learner and $p_d(\dot{\mathbf{X}}; \boldsymbol{\theta}_d)$ for the data generator, that governs the players' strategy selection. Players are allowed to select the parameterization $\boldsymbol{\theta}_l$ and $\boldsymbol{\theta}_d$ for the respective distributions. For any choice of $\boldsymbol{\theta}_l$, the l -player plays a strategy \mathbf{w} sampled from $p_l(\cdot; \boldsymbol{\theta}_l)$. Similarly, for any choice of $\boldsymbol{\theta}_d$, the d -player plays a strategy $\dot{\mathbf{X}}$ sampled from $p_d(\cdot; \boldsymbol{\theta}_d)$. If the players adhere to such rules, we obtain a randomized prediction game.

A randomized prediction game is a noncooperative game between a learner (l -player) and data generator (d -player) that has the following components.

- 1) An underlying prediction game with cost functions $c_{l/d}(\mathbf{w}, \dot{\mathbf{X}})$ as defined in (1) and (2).
- 2) two parametrized probability distributions $p_{l/d}(\cdot; \boldsymbol{\theta}_{l/d})$ with parameters in $\Theta_{l/d}$.

¹We adopt here the same terminology used in [11].

3) $\Theta_{l/d}$ are nonempty, compact, and convex subsets of a finite-dimensional metric space $\mathbb{R}^{S_{l/d}}$.

The sets of parameters $\Theta_{l/d}$ are the pure strategy sets (also known as action spaces) for the l -player and d -player, respectively. The cost functions of the two players, which quantify the cost that each player incurs when a strategy profile $(\theta_l, \theta_d) \in \Theta_l \times \Theta_d$ is played, coincide with the expected costs, denoted by $\bar{c}_{l/d}(\theta_l, \theta_d)$, that the two players have in the underlying prediction game if strategies are sampled from $p_l(\cdot; \theta_l)$ and $p_d(\cdot; \theta_d)$, according to the expected cost-minimizing hypothesis

$$\bar{c}_l(\theta_l, \theta_d) = \mathbb{E}_{\mathbf{w} \sim p_l(\cdot; \theta_l)} [c_l(\mathbf{w}, \dot{\mathbf{X}})] \quad (3)$$

$$\bar{c}_d(\theta_l, \theta_d) = \mathbb{E}_{\mathbf{w} \sim p_l(\cdot; \theta_l)} [c_d(\mathbf{w}, \dot{\mathbf{X}})] \quad (4)$$

where $\mathbb{E}[\cdot]$ denotes the expectation operator. We assume $\bar{c}_{l/d}$ to be well-defined functions, i.e., the expectations to be finite for any $(\theta_l, \theta_d) \in \Theta_l \times \Theta_d$. To avoid confusion between $c_{l/d}$ and $\bar{c}_{l/d}$, in the following, we will refer them, respectively, as cost functions, and expected cost functions.

By adhering to a noncooperative setting, the two players involved in the prediction game are not allowed to communicate and they play their strategies simultaneously. Each player has complete information of the game setting by knowing the expected cost function and strategy set of either players. Under rationality assumption, each player's interest is to achieve the greatest personal advantage, i.e., to incur the lowest possible cost. Accordingly, the players are prone to play a Nash equilibrium, which in the context of our randomized prediction game is a strategy profile $(\theta_l^*, \theta_d^*) \in \Theta_l \times \Theta_d$, such that no player is interested in changing his/her own playing strategy. In formal terms, this yields

$$\theta_l^* \in \arg \min_{\theta_l \in \Theta_l} \bar{c}_l(\theta_l, \theta_d^*), \quad \theta_d^* \in \arg \min_{\theta_d \in \Theta_d} \bar{c}_d(\theta_l^*, \theta_d). \quad (5)$$

B. Existence of a Nash Equilibrium

The existence of a Nash equilibrium of a randomized prediction game is not granted in general. A sufficient condition is given in the following.

Theorem 1 (Existence): A randomized prediction game admits at least one Nash equilibrium if it has the following.

- 1) $\bar{c}_{l/d}$ are continuous in $\Theta_l \times \Theta_d$.
- 2) $\bar{c}_l(\cdot, \theta_d)$ is quasi-convex in Θ_l for any $\theta_d \in \Theta_d$.
- 3) $\bar{c}_d(\theta_l, \cdot)$ is quasi-convex in Θ_d for any $\theta_l \in \Theta_l$.

Proof: The result follows directly from the Debreu–Glicksberg–Fan theorem [29]. ■

C. Uniqueness of a Nash Equilibrium

In addition to the existence of a Nash equilibrium, it is of interest to investigate if the equilibrium is unique. However, determining tight conditions that guarantee the uniqueness of the Nash equilibrium for any randomized prediction game is challenging; in particular, due to the additional dependence on a probability distribution for the learner and the data generator.

We will make use of a classical result due to Rosen [30] to formulate sufficient conditions for the uniqueness of the Nash

equilibrium of randomized prediction games in terms of the so-called pseudogradient of the game, defined as

$$\bar{\mathbf{g}}_{\mathbf{r}} = \begin{bmatrix} r_l \nabla_{\theta_l} \bar{c}_l \\ r_d \nabla_{\theta_d} \bar{c}_d \end{bmatrix} \quad (6)$$

with any fixed vector $\mathbf{r} = [r_l, r_d]^\top \geq \mathbf{0}$. In particular, a randomized prediction game admits a unique Nash equilibrium if the following assumption is verified.

Assumption 1:

- 1) $\bar{c}_{l/d}$ are twice differentiable in $\Theta_l \times \Theta_d$.
- 2) $\bar{c}_l(\cdot, \theta_d)$ is convex in Θ_l for any $\theta_d \in \Theta_d$.
- 3) $\bar{c}_d(\theta_l, \cdot)$ is convex in Θ_d for any $\theta_l \in \Theta_l$

and $\bar{\mathbf{g}}_{\mathbf{r}}$ is strictly monotone for some fixed $\mathbf{r} > \mathbf{0}$, that is,

$$[\bar{\mathbf{g}}_{\mathbf{r}}(\theta_l, \theta_d) - \bar{\mathbf{g}}_{\mathbf{r}}(\theta'_l, \theta'_d)]^\top \begin{bmatrix} \theta_l - \theta'_l \\ \theta_d - \theta'_d \end{bmatrix} > 0$$

for any distinct strategy profiles $(\theta_l, \theta_d), (\theta'_l, \theta'_d) \in \Theta_l \times \Theta_d$.²

Rosen [30] also provides a useful sufficient condition that guarantees a strictly monotone pseudogradient. This requires the Jacobian of the pseudogradient, also known as pseudo-Jacobian, given by

$$\bar{\mathbf{J}}_{\mathbf{r}} = \begin{bmatrix} r_l \nabla_{\theta_l}^2 \bar{c}_l & r_l \nabla_{\theta_l \theta_d}^2 \bar{c}_l \\ r_d \nabla_{\theta_d}^2 \bar{c}_d & r_d \nabla_{\theta_d \theta_l}^2 \bar{c}_d \end{bmatrix} \quad (7)$$

to be positive definite.

Theorem 2: A randomized prediction game admits a unique Nash equilibrium if Assumption 1 holds, and the pseudo-Jacobian $\bar{\mathbf{J}}_{\mathbf{r}}(\theta_l, \theta_d)$ is positive definite for all $(\theta_l, \theta_d) \in \Theta_l \times \Theta_d$ and some fixed $\mathbf{r} > \mathbf{0}$.

Proof: Under Assumption 1, the positive definiteness of $\bar{\mathbf{J}}_{\mathbf{r}}$ for all strategy profiles and some fixed vector $\mathbf{r} > \mathbf{0}$ implies the strict monotonicity of $\bar{\mathbf{g}}_{\mathbf{r}}$, which in turn implies the uniqueness of the Nash equilibrium [30, Th. 6]. ■

In the following, we provide sufficient conditions that ensure the positive definiteness of the pseudo-Jacobian and the uniqueness of the Nash equilibrium via Theorem. 2. To this end, we decompose $\bar{c}_{l/d}(\theta_l, \theta_d)$ as:

$$\begin{aligned} \bar{c}_l(\theta_l, \theta_d) &= \rho_l \bar{\Omega}_l(\theta_l) + \bar{L}_l(\theta_l, \theta_d) \\ \bar{c}_d(\theta_l, \theta_d) &= \rho_d \bar{\Omega}_d(\theta_d) + \bar{L}_d(\theta_l, \theta_d) \end{aligned} \quad (8)$$

where $\bar{\Omega}_{l/d}$ and $\bar{L}_{l/d}$ are the expected regularization and loss terms given by

$$\begin{aligned} \bar{\Omega}_l(\theta_l) &= \mathbb{E}_{\mathbf{w} \sim p_l(\cdot; \theta_l)} [\Omega_l(\mathbf{w})] \\ \bar{\Omega}_d(\theta_d) &= \mathbb{E}_{\dot{\mathbf{X}} \sim p_d(\cdot; \theta_d)} [\Omega_d(\dot{\mathbf{X}})] \\ \bar{L}_l(\theta_l, \theta_d) &= \mathbb{E}_{\mathbf{w} \sim p_l(\cdot; \theta_l)} \left[\sum_{\dot{\mathbf{X}} \sim p_d(\cdot; \theta_d)}^n \ell_l(\mathbf{w}^\top \mathbf{x}_i, y_i) \right] \\ \bar{L}_d(\theta_l, \theta_d) &= \mathbb{E}_{\mathbf{w} \sim p_l(\cdot; \theta_l)} \left[\sum_{\dot{\mathbf{X}} \sim p_d(\cdot; \theta_d)}^n \ell_d(\mathbf{w}^\top \mathbf{x}_i, y_i) \right]. \end{aligned}$$

Moreover, we require the following convexity and differentiability conditions on $\bar{\Omega}_{l/d}$ and $\bar{L}_{l/d}$.

²Assumption 1-1) could be relaxed to continuously differentiable.

Assumption 2:

- 1) $\bar{\Omega}_{l/d}$ is strongly convex and twice continuously differentiable in $\Theta_{l/d}$.
- 2) $\bar{L}_l(\cdot, \theta_d)$ is convex and twice continuously differentiable in Θ_l for all $\theta_d \in \Theta_d$.
- 3) $\bar{L}_d(\theta_l, \cdot)$ is convex and twice continuously differentiable in Θ_d for all $\theta_l \in \Theta_l$.

Finally, we introduce some quantities that are used in the subsequent lemma, which gives sufficient conditions for the positive definiteness of the pseudo-Jacobian

$$\begin{aligned}\lambda_l^\Omega &= \inf_{\theta_l \in \Theta_l} \lambda_{\min}[\nabla_{\theta_l}^2 \bar{\Omega}_l(\theta_l)] \\ \lambda_d^\Omega &= \inf_{\theta_d \in \Theta_d} \lambda_{\min}[\nabla_{\theta_d}^2 \bar{\Omega}_d(\theta_d)] \\ \lambda_l^L &= \inf_{(\theta_l, \theta_d) \in \Theta_l \times \Theta_d} \lambda_{\min}[\nabla_{\theta_l}^2 \bar{L}_l(\theta_l, \theta_d)] \\ \lambda_d^L &= \inf_{(\theta_l, \theta_d) \in \Theta_l \times \Theta_d} \lambda_{\min}[\nabla_{\theta_d}^2 \bar{L}_d(\theta_l, \theta_d)] \\ \tau &= \sup_{(\theta_l, \theta_d) \in \Theta_l \times \Theta_d} \lambda_{\max}[R(\theta_l, \theta_d)R(\theta_l, \theta_d)^\top]\end{aligned}$$

where $R(\theta_l, \theta_d) = 1/2[\nabla_{\theta_l}^2 \bar{L}_l(\theta_l, \theta_d)^\top + \nabla_{\theta_d}^2 \bar{L}_d(\theta_l, \theta_d)]$ and $\lambda_{\max/\min}$ give the maximum/minimum eigenvalue of the matrix in input. The quantities listed above are finite and positive if Assumption 2 holds, given the compactness of $\Theta_{l/d}$.

Lemma 1: If Assumption 2 holds and

$$(\rho_l \lambda_l^\Omega + \lambda_l^L)(\rho_d \lambda_d^\Omega + \lambda_d^L) > \tau$$

then the pseudo-Jacobian $\bar{J}_r(\theta_l, \theta_d)$ is positive definite for all $(\theta_l, \theta_d) \in \Theta_l \times \Theta_d$ by taking $\mathbf{r} = (1, 1)^\top$.

Proof: The pseudo-Jacobian in (7) can be written as follows given the decomposition of $\bar{c}_{l/d}$ in (8):

$$\mathbf{J}_r = \begin{bmatrix} \rho_l \nabla_{\theta_l}^2 \bar{\Omega}_l + \nabla_{\theta_l}^2 \bar{L}_l & \nabla_{\theta_l \theta_d}^2 \bar{L}_l \\ \nabla_{\theta_d \theta_l}^2 \bar{L}_d & \rho_d \nabla_{\theta_d}^2 \bar{\Omega}_d + \nabla_{\theta_d \theta_l}^2 \bar{L}_d \end{bmatrix}$$

where we omitted the arguments of $\bar{\Omega}_{l/d}$ and $\bar{L}_{l/d}$ for notational convenience. Let us denote by \mathbf{J}_r^{ll} , \mathbf{J}_r^{ld} , \mathbf{J}_r^{dl} , and \mathbf{J}_r^{dd} the four matrices composing \mathbf{J}_r (in top-down, left-right order).

Consider the following matrix:

$$\mathbf{H} = \begin{bmatrix} \mathbf{H}^{ll} & \mathbf{H}^{ld} \\ \mathbf{H}^{dl} & \mathbf{H}^{dd} \end{bmatrix} = \begin{bmatrix} \rho_l \lambda_l^\Omega + \lambda_l^L & R(\theta_l, \theta_d)^\top \\ R(\theta_l, \theta_d) & \rho_d \lambda_d^\Omega + \lambda_d^L \end{bmatrix}.$$

Then, we have for all $\mathbf{t} = (\mathbf{t}_l^\top, \mathbf{t}_d^\top)^\top \neq \mathbf{0}$

$$\begin{aligned}\mathbf{t}^\top \mathbf{J}_r \mathbf{t} &= \mathbf{t}^\top \frac{\mathbf{J}_r + \mathbf{J}_r^\top}{2} \mathbf{t}^\top \\ &= \underbrace{\mathbf{t}_l^\top \mathbf{J}_r^{ll} \mathbf{t}_l}_{\geq \mathbf{t}_l^\top \mathbf{H}^{ll} \mathbf{t}_l} + \underbrace{\mathbf{t}_d^\top \mathbf{J}_r^{dd} \mathbf{t}_d}_{\geq \mathbf{t}_d^\top \mathbf{H}^{dd} \mathbf{t}_d} + \underbrace{\mathbf{t}_l^\top \frac{\mathbf{J}_r^{ld} + \mathbf{J}_r^{dl\top}}{2} \mathbf{t}_d}_{\mathbf{H}^{ld} + \mathbf{H}^{dl\top}} \geq \mathbf{t}^\top \mathbf{H} \mathbf{t}\end{aligned}$$

where the underbraced relations follow from the definitions of $\lambda_{l/d}^\Omega$, $\lambda_{l/d}^L$, and R . Accordingly, the positive definiteness of \mathbf{J}_r can be derived from the positive definiteness of matrix \mathbf{H} . To prove the latter, we will show that all roots of the characteristic polynomial $\det(\mathbf{H} - \lambda \mathbf{I})$ of \mathbf{H} are positive. By properties

of the determinant,³ we have

$$\det(\mathbf{H} - \lambda \mathbf{I}) = \det\left((\rho_l \lambda_l^\Omega + \lambda_l^L - \lambda) \mathbf{I} \cdot \det\left((\rho_d \lambda_d^\Omega + \lambda_d^L - \lambda) \mathbf{I} - \frac{\mathbf{S}}{\rho_l \lambda_l^\Omega + \lambda_l^L - \lambda}\right)\right)$$

where \mathbf{S} is a diagonal matrix with the eigenvalues of $R(\theta_l, \theta_d)R(\theta_l, \theta_d)^\top$. The roots of the first determinant term are all equal to $\rho_l \lambda_l^\Omega + \lambda_l^L$, which is positive because $\rho_l > 0$ by construction and $\lambda_l^\Omega > 0$ follows from the strong convexity of $\bar{\Omega}_l$ in Assumption 2-1). As for the second determinant term, take the i th diagonal element \mathbf{S}_{ii} of \mathbf{S} . Then, two roots are the solution of the following quadratic equation:

$$\lambda^2 - \lambda(a + b) + ab - \mathbf{S}_{ii} = 0$$

which are given by

$$\lambda_{1,2}^{(i)} = a + b \pm \sqrt{(a - b)^2 + 4\mathbf{S}_{ii}}$$

where $a = \rho_l \lambda_l^\Omega + \lambda_l^L$ and $b = \rho_d \lambda_d^\Omega + \lambda_d^L$. Among the two, $\lambda_2^{(i)}$ (the one with the minus) is the smallest one, which is strictly positive if

$$ab = (\rho_l \lambda_l^\Omega + \lambda_l^L)(\rho_d \lambda_d^\Omega + \lambda_d^L) > \mathbf{S}_{ii}.$$

Since the condition has to hold for any choice of the eigenvalue \mathbf{S}_{ii} in the right-hand side of the inequality, we take the maximum one $\max_i \mathbf{S}_{ii}$, which coincides with $\lambda_{\max}(R(\theta_l, \theta_d)R(\theta_l, \theta_d)^\top)$. We further maximize the latter quantity with respect to $(\theta_l, \theta_d) \in \Theta_l \times \Theta_d$, as we want the result to hold for any parameterization. Therefrom, we recover the variable τ and the condition $(\rho_l \lambda_l^\Omega + \lambda_l^L)(\rho_d \lambda_d^\Omega + \lambda_d^L) > \tau$, which guarantees that all roots of the characteristic polynomial of \mathbf{H} are strictly positive for any choice of $(\theta_l, \theta_d) \in \Theta_l \times \Theta_d$, and hence, \bar{J}_r is positive definite over $\Theta_l \times \Theta_d$. ■

In addition to Lemma 1, we provide in the Supplementary Material alternative (stronger) sufficient conditions, which generalize the ones given in [11].

D. Finding a Nash Equilibrium

From the computational perspective, we can find a Nash equilibrium in our game by exploiting algorithms similar to the ones adopted for static prediction games [11]. In particular, we consider a modified extragradient descent algorithm [11], [31], [32] that finds a solution to the following variational inequality problem, provided that \bar{g}_r is continuous and monotone:

$$\bar{g}_r(\theta_l^*, \theta_d^*)^\top (\theta - \theta^*) \geq 0 \quad \forall (\theta_l, \theta_d) \in \Theta_l \times \Theta_d \quad (9)$$

where $\theta = [\theta_l^\top, \theta_d^\top]^\top$ and similarly for θ^* . Any solution θ^* to this problem can be shown to correspond bijectively to a Nash equilibrium of a game having \bar{g}_r as pseudogradient [11], [32].

If Theorem 1 holds, the pseudo-Jacobian \bar{J}_r can be shown to be positive semidefinite, and \bar{g}_r is thus continuous and

³ $\det\left[\frac{a}{b} \mathbf{I} \mathbf{B}^\top\right] = \det(a \mathbf{I}) \det(d \mathbf{I} - \frac{1}{a} \mathbf{B} \mathbf{B}^\top)$ and if $\mathbf{U} \mathbf{S} \mathbf{U}^\top$ is the eigendecomposition of $\mathbf{B} \mathbf{B}^\top$, then the latter determinant becomes $\det(\mathbf{U}(d \mathbf{I} - \frac{1}{a} \mathbf{S}) \mathbf{U}^\top) = \det(d \mathbf{I} - \frac{1}{a} \mathbf{S})$

Algorithm 1 Extragradient Descent (Adapted From [11])

Input: Cost functions $\bar{c}_{l/d}$; parameter spaces Θ_l, Θ_d ; a small positive constant ϵ .

Output: The optimal parameters θ_l, θ_d .

- 1: Randomly select $\theta^{(0)} = (\theta_l^{(0)}, \theta_d^{(0)}) \in \Theta_l \times \Theta_d$.
- 2: Set iteration count $k = 0$, and select $\sigma, \beta \in (0, 1)$.
- 3: Set $\mathbf{r} = (r_l, r_d)^\top = (1, \rho_l/\rho_d)^\top$.
- 4: **repeat**
- 5: Set $\mathbf{d}^{(k)} = \Pi_{\Theta_l \times \Theta_d} \left(\theta^{(k)} - \bar{\mathbf{g}}_{\mathbf{r}} \left(\theta_l^{(k)}, \theta_d^{(k)} \right) \right) - \theta^{(k)}$.
- 6: Find maximum step size $t^{(k)} \in \{\beta^p | p \in \mathbb{N}\}$ s.t.

$$-\bar{\mathbf{g}}_{\mathbf{r}} \left(\bar{\theta}_l^{(k)}, \bar{\theta}_d^{(k)} \right)^\top \mathbf{d}^{(k)} \geq \sigma \left(\|\mathbf{d}^{(k)}\|_2^2 \right),$$

where $\bar{\theta}^{(k)} = \theta^{(k)} + t^{(k)} \mathbf{d}^{(k)}$.

- 7: Set $\eta^{(k)} = -\frac{t^{(k)}}{\|\bar{\mathbf{g}}_{\mathbf{r}}(\bar{\theta}_l^{(k)}, \bar{\theta}_d^{(k)})\|_2} \bar{\mathbf{g}}_{\mathbf{r}} \left(\bar{\theta}_l^{(k)}, \bar{\theta}_d^{(k)} \right)^\top \mathbf{d}^{(k)}$.
- 8: Set $\theta^{(k+1)} = \Pi_{\Theta_l \times \Theta_d} \left(\theta^{(k)} - \eta^{(k)} \bar{\mathbf{g}}_{\mathbf{r}} \left(\bar{\theta}_l^{(k)}, \bar{\theta}_d^{(k)} \right) \right)$.
- 9: Set $k = k + 1$.
- 10: **until** $\|\theta^{(k)} - \theta^{(k-1)}\|_2 \leq \epsilon$
- 11: **return** $\theta_l = \theta_l^{(k)}, \theta_d = \theta_d^{(k)}$

monotone. Hence, the variational inequality can be solved by the modified extragradient descent algorithm given as Algorithm 1, which is guaranteed to converge to a Nash equilibrium point [31], [33]. The algorithm generates a sequence of feasible points whose distance from the equilibrium solution is monotonically decreased. It exploits a projection operator $\Pi_{\Theta_l \times \Theta_d}(\theta)$ to map the input vector θ onto the closest admissible point in $\Theta_l \times \Theta_d$, and a simple line-search algorithm to find the maximum step t on the descent direction \mathbf{d} .⁴

In Section III, we apply our randomized prediction game to the case of linear SVM learners, and compute the corresponding pseudogradient, as required in Algorithm 1.

III. RANDOMIZED PREDICTION GAMES FOR SUPPORT VECTOR MACHINES

In this section, we consider a randomized prediction game involving a linear SVM learner [34] and Gaussian distributions as the underlying probabilities $p_{l/d}$.

Learner: The decision function of the learner is of the type $f(\mathbf{x}; \mathbf{w}) = \mathbf{w}^\top \phi(\mathbf{x})$ where the feature map is given by $\phi(\mathbf{x}) = [\mathbf{x}^\top \ 1]^\top$. For convenience, we consider a decomposition of \mathbf{w} into $[\tilde{\mathbf{w}}^\top \ b]^\top$, where $\tilde{\mathbf{w}} \in \mathbb{R}^{m-1}$ and $b \in \mathbb{R}$. Hence, the decision function can also be written as $f(\mathbf{x}; \mathbf{w}) = \tilde{\mathbf{w}}^\top \mathbf{x} + b$. Accordingly, the input space \mathcal{X} is a $(m-1)$ -dimensional vector space, i.e., $\mathcal{X} \subseteq \mathbb{R}^{m-1}$. The distribution p_l for the learner is assumed to be Gaussian. In order to guarantee the theoretical existence of the Nash equilibrium through Theorem 1, we assume the parameters of the Gaussian distribution to be bounded. For the sake of clarity, we use in this section axis-aligned Gaussians (i.e., with diagonal covariance matrices) for our analysis, even though general covariances could be adopted as well. Under these assumptions, we define the

strategy set for the learner as $\Theta_l = \{(\mu_{\mathbf{w}}, \sigma_{\mathbf{w}}) \in \mathbb{R}^m \times \mathbb{R}_+^m\} \cap \mathcal{B}_l$, where $\mathcal{B}_l \subset \mathbb{R}^m \times \mathbb{R}_+^m$ is an application-dependent nonempty, convex, bounded set, restricting the set of feasible parameters. The parameter vectors $\mu_{\mathbf{w}}$ and $\sigma_{\mathbf{w}}$ encode the mean and standard deviation of the axis-aligned Gaussian distributions. The loss function ℓ_l of the learner corresponds to the hinge loss of the SVM, i.e., $\ell_l(z, y) = [1 - zy]_+$ with $[z]_+ = \max(0, z)$, while the strategy penalization term $\Omega_l(\mathbf{w})$ is the squared Euclidean norm of $\tilde{\mathbf{w}}$. As a result, the cost function c_l corresponds to the C-SVM objective function, and it is convex in \mathbf{w}

$$c_l(\mathbf{w}, \mathbf{X}) = \frac{\rho_l}{2} \|\tilde{\mathbf{w}}\|^2 + \sum_{i=1}^n [1 - y_i(\tilde{\mathbf{w}}^\top \mathbf{x}_i + b)]_+. \quad (10)$$

Data Generator: For convenience, we consider \mathbf{X} rather than $\dot{\mathbf{X}}$ as the quantity undergoing the randomization. This comes without loss of generality, because there is a one-to-one correspondence between $\dot{\mathbf{x}}_i$ and \mathbf{x}_i if we consider the linear feature map $\dot{\mathbf{x}}_i = \phi(\mathbf{x}_i) = [\mathbf{x}_i^\top, 1]^\top$. Moreover, we assume that samples \mathbf{x}_i can be perturbed independently. Accordingly, the distribution p_d for the data generator factorizes as $p_d(\mathbf{X}; \theta_d) = \prod_{i=1}^n p_d(\mathbf{x}_i; \theta_d^{(i)})$, where $\theta_d = (\theta_d^{(1)}, \dots, \theta_d^{(n)})$. We consider $p_d(\mathbf{x}_i; \theta_d^{(i)})$ to be a k -variate axis-aligned Gaussian distribution with bounded mean and standard deviation given by $\theta_d^{(i)} = (\mu_{\mathbf{x}_i}, \sigma_{\mathbf{x}_i})$. In summary, the strategy set adopted for the data generator is given by $\Theta_d = \prod_{i=1}^n \Theta_d^{(i)}$, where $\Theta_d^{(i)} = \{(\mu_{\mathbf{x}_i}, \sigma_{\mathbf{x}_i}) \in \mathbb{R}^k \times \mathbb{R}_+^k\} \cap \mathcal{B}_d$. Here, $\mathcal{B}_d \subset \mathbb{R}^k \times \mathbb{R}_+^k$ is a nonempty, convex, bounded set. The loss function ℓ_d of the data generator is the hinge loss under wrong labeling, i.e., $\ell_d(z, y) = [1 + zy]_+$. In this way, the data generator is penalized if the learner correctly classifies a sample point. Finally, the strategy penalization function Ω_d is the squared Euclidean distance of the perturbed samples in \mathbf{X} from the ones in the original training set $\hat{\mathcal{D}}$, i.e., $\Omega_d(\mathbf{X}) = \sum_{i=1}^n \|\mathbf{x}_i - \hat{\mathbf{x}}_i\|^2$. The resulting cost function c_d is convex in \mathbf{X}

$$c_d(\mathbf{w}, \mathbf{X}) = \frac{\rho_d}{2} \sum_{i=1}^n \|\mathbf{x}_i - \hat{\mathbf{x}}_i\|^2 + \sum_{i=1}^n [1 + y_i(\tilde{\mathbf{w}}^\top \mathbf{x}_i + b)]_+. \quad (11)$$

Existence of a Nash Equilibrium: The proposed randomized prediction game for the SVM learner admits at least one Nash equilibrium. This can be proven by means of Theorem 1. Indeed, the required continuity of $\bar{c}_{l/d}$ holds and, as for the quasi-convexity conditions, we can rewrite (3) as follows by exploiting the fact that p_l is a Gaussian distribution with mean $\mu_{\mathbf{w}}$ and standard deviation $\sigma_{\mathbf{w}}$:

$$\bar{c}_l(\theta_l, \theta_d) = \mathbb{E}_{\substack{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, I) \\ \mathbf{X} \sim p_d(\cdot; \theta_d)}} [c_l(\mu_{\mathbf{w}} + D(\sigma_{\mathbf{w}})\mathbf{z}, \mathbf{X})] \quad (12)$$

where $\mathcal{N}(\mathbf{0}, I)$ is a m -dimensional standard normal distribution and $D(\sigma_{\mathbf{w}})$ is a diagonal matrix having $\sigma_{\mathbf{w}}$ on the diagonal. Since c_l is convex in its first argument and convexity is preserved under addition of convex functions, positive rescaling, and composition with linear functions, we have that

⁴We refer the reader to [11], [31], [32] (and references therein) for proofs that derive conditions for which \mathbf{d} is effectively a descent direction.

\bar{c}_l is convex (and thus quasi-convex) in $\theta_l = (\mu_w, \sigma_w)$. As for the quasi-convexity condition of the data generator's cost, we can exploit the separability of c_d to rewrite (4) as follows:

$$\bar{c}_d(\theta_l, \theta_d) = \sum_{i=1}^n \mathbb{E}_{\substack{\mathbf{w} \sim p_l(\cdot; \theta_l) \\ \mathbf{z} \sim \mathcal{N}(\mathbf{0}, I)}} [c_d^{(i)}(\mathbf{w}, \mu_{x_i} + D(\sigma_{x_i})\mathbf{z})]$$

where

$$c_d^{(i)}(\mathbf{w}, \mathbf{x}) = \frac{\rho_d}{2} \|\mathbf{x} - \hat{\mathbf{x}}_i\|^2 + [1 + y_i(\mathbf{w}^\top \hat{\mathbf{x}}_i + b)]_+.$$

Since $c_d^{(i)}$ is convex in its second argument, by following the same reasoning used to show the quasi-convexity of the learner's expected cost, we have that each expectation in \bar{c}_d is convex in $\theta_d^{(i)} = (\mu_{x_i}, \sigma_{x_i})$, $1 \leq i \leq n$. As a consequence, \bar{c}_d is convex, and hence, quasi-convex in θ_d is the sum of convex functions.

Uniqueness of a Nash Equilibrium: We have shown that $\bar{c}_l(\cdot, \theta_d)$ and $\bar{c}_d(\theta_l, \cdot)$ are convex as required by Assumption 1-2) and 1-3). In particular, we have that the single expected regularization terms $\bar{\Omega}_{l/d}(\cdot)$ and loss terms $\bar{L}_l(\cdot, \theta_d)$, $\bar{L}_d(\theta_l, \cdot)$ are convex as well. Moreover, they are twice continuously differentiable by having Gaussian distributions for $p_{l/d}$. It is then sufficient to have $\bar{\Omega}_{l/d}$, which are strongly convex to prove the uniqueness of the Nash equilibrium via Lemma 1. While it is easy to see that $\bar{\Omega}_d$ is strongly convex, we have that $\bar{\Omega}_l$ is not strongly convex with respect to b due to the presence of an unregularized bias term b in the learner. The problem derives from the fact that the SVM itself may not have a unique solution when the bias term is present and nonregularized (see [35], [36] for a characterization of the degenerate cases). As a result, the proposed game is not guaranteed to have a unique Nash equilibrium in its actual form. On the other hand, a unique Nash equilibrium may be obtained by either considering an unbiased SVM, i.e., by setting $b = 0$ as in [11], or a regularized bias term, e.g., by adding $\frac{\epsilon}{2}b^2$ to the learner's objective function with $\epsilon > 0$. In both cases, all conditions that ensure the uniqueness of the Nash equilibrium via Theorem 2 and Lemma 1 would be satisfied, under proper choices of $p_{l/d}$.

It is worth noting however that the necessary and sufficient conditions under which a biased (nonregularized) SVM has no unique solution are quite restricted [35], [36]. For this reason, we believe that uniqueness of the Nash equilibrium could also be proven for the biased SVM under mild assumptions. However, this requires considerable effort in trying to relax the sufficiency conditions of Rosen [30], which is beyond the scope of this paper. We thus leave this challenge to future investigations. Moreover, we believe that enforcing a unique Nash equilibrium in our game by making the original SVM formulation strictly convex may lead to worse results, similar to exploiting convex approximations to solve originally nonconvex problems in machine learning [37], [38]. For the above reasons, in this paper, we choose to retain the original SVM formulation for the learner, by sacrificing the uniqueness of the Nash Equilibrium. We nevertheless provide in Section V a discussion of why having a unique Nash Equilibrium is not so important in practice for our game, and we empirically show in Section VI that our approach can anyway achieve

competitive performances with respect to other state-of-the-art approaches.

The rest of this section is devoted to showing how to compute the pseudogradient (6) by providing explicit formulas for $\nabla_{\theta_l} \bar{c}_l$ and $\nabla_{\theta_d} \bar{c}_d$.

A. Gradient of the Learner's Cost

In this section, we focus on computing the gradient $\nabla_{\theta_l} \bar{c}_l(\theta_l, \theta_d)$, where \bar{c}_l is defined as in (10). By properties of expectation and since \mathbf{w} follows an axis-aligned Gaussian distribution with mean μ_w and standard deviation σ_w , we can reduce the cost of the learner to:

$$\bar{c}_l(\theta_l, \theta_d) = \frac{\rho_l}{2} (\|\mu_{\tilde{\mathbf{w}}}\|^2 + \|\sigma_{\tilde{\mathbf{w}}}\|^2) + \sum_{i=1}^n \mathbb{E}_{\substack{\mathbf{w} \sim p_l(\cdot; \theta_l) \\ \mathbf{x}_i \sim p_d(\cdot; \theta_d^{(i)})}} [[1 - y_i(\tilde{\mathbf{w}}^\top \mathbf{x}_i + b)]_+] \quad (13)$$

where we are assuming the following decompositions for the mean $\mu_w = [\mu_{\tilde{\mathbf{w}}}^\top \mu_b]^\top$ and standard deviation $\sigma_w = [\sigma_{\tilde{\mathbf{w}}}^\top \sigma_b]^\top$. The hard part for the minimization is the term in the expectation, which cannot be expressed to our knowledge in a closed-form function of the Gaussian's parameters. We thus resort to a central-limit-theorem-like approximation, by regarding $s_i = 1 - y_i(\tilde{\mathbf{w}}^\top \mathbf{x}_i + b)$ as a Gaussian-distributed variable with mean μ_{s_i} and standard deviation σ_{s_i} , i.e., $s_i \sim \mathcal{N}(\mu_{s_i}, \sigma_{s_i})$. In general, s_i does not follow a Gaussian distribution, since the product of two normal deviates is not normally distributed. However, if the number of features k is large, the approximation becomes reasonable. Under this assumption, we can rewrite the expectation as follows:

$$\begin{aligned} & \mathbb{E}_{\substack{\mathbf{w} \sim p_l(\cdot; \theta_l) \\ \mathbf{x}_i \sim p_d(\cdot; \theta_d^{(i)})}} [[1 - y_i(\tilde{\mathbf{w}}^\top \mathbf{x}_i + b)]_+] \\ &= \mathbb{E}_{s_i \sim \mathcal{N}(\mu_{s_i}, \sigma_{s_i})} [[s_i]_+]. \end{aligned} \quad (14)$$

The mean and the variance of the Gaussian distribution in the right-hand side of (14) are, respectively, given by

$$\begin{aligned} \mu_{s_i} &= \mathbb{E}_{\substack{\mathbf{w} \sim p_l(\cdot; \theta_l) \\ \mathbf{x}_i \sim p_d(\cdot; \theta_d^{(i)})}} [1 - y_i(\tilde{\mathbf{w}}^\top \mathbf{x}_i + b)] \\ &= 1 - y_i(\mu_{\tilde{\mathbf{w}}}^\top \mu_{x_i} + \mu_b) \end{aligned} \quad (15)$$

$$\begin{aligned} \sigma_{s_i}^2 &= \mathbb{V}_{\substack{\mathbf{w} \sim p_l(\cdot; \theta_l) \\ \mathbf{x}_i \sim p_d(\cdot; \theta_d^{(i)})}} [1 - y_i(\tilde{\mathbf{w}}^\top \mathbf{x}_i + b)] \\ &= \sigma_{\tilde{\mathbf{w}}}^2 (\sigma_{x_i}^2 + \mu_{x_i}^2) + \mu_{\tilde{\mathbf{w}}}^2 \sigma_{x_i}^2 + \sigma_b^2 \end{aligned} \quad (16)$$

where \mathbb{V} is the variance operator, and we assume that squaring a vector corresponds to squaring each single component.

The expectation in (14) can be transformed after simple manipulations into the following function involving the Gauss error function (integral function of the standard normal distribution) denoted by $\text{erf}()$:

$$\begin{aligned} h(\mu_{s_i}, \sigma_{s_i}) &= \frac{\sigma_{s_i}}{\sqrt{2\pi}} \exp\left(-\frac{\mu_{s_i}^2}{2\sigma_{s_i}^2}\right) \\ &+ \frac{\mu_{s_i}}{2} \left[1 - \text{erf}\left(-\frac{\mu_{s_i}}{\sqrt{2}\sigma_{s_i}}\right)\right]. \end{aligned} \quad (17)$$

The learner's cost in (13) can thus be approximated as

$$\bar{c}_l(\theta_l, \theta_d) \approx L_l(\mu_w, \sigma_w) = \frac{\rho_l}{2} (\|\mu_{\tilde{w}}\|^2 + \|\sigma_{\tilde{w}}\|^2) + \sum_{i=1}^n h(\mu_{s_i}(\theta_l), \sigma_{s_i}(\theta_l)). \quad (18)$$

We can now approximate the gradient $\nabla_{\theta_l} \bar{c}_l$ in terms of $\nabla_{\theta_l} L_l$. In the following, we denote the Hadamard (also known as entrywise) product between any two vectors \mathbf{a} and \mathbf{b} as $\mathbf{a} \circ \mathbf{b}$, and we assume any scalar-by-vector derivative to be a column vector. The gradients of interest are given as

$$\frac{\partial L_l}{\partial \mu_w} = \rho_l \begin{bmatrix} \mu_{\tilde{w}} \\ 0 \end{bmatrix} + \sum_{i=1}^n \left(\frac{\partial h}{\partial \mu_{s_i}} \frac{\partial \mu_{s_i}}{\partial \mu_w} + \frac{\partial h}{\partial \sigma_{s_i}^2} \frac{\partial \sigma_{s_i}^2}{\partial \mu_w} \right) \quad (19)$$

$$\frac{\partial L_l}{\partial \sigma_w} = \rho_l \begin{bmatrix} \sigma_{\tilde{w}} \\ 0 \end{bmatrix} + \sum_{i=1}^n \left(\frac{\partial h}{\partial \mu_{s_i}} \frac{\partial \mu_{s_i}}{\partial \sigma_w} + \frac{\partial h}{\partial \sigma_{s_i}^2} \frac{\partial \sigma_{s_i}^2}{\partial \sigma_w} \right) \quad (20)$$

where it is not difficult to show that

$$\frac{\partial h}{\partial \mu_{s_i}} = \frac{1}{2} \left[1 - \operatorname{erf} \left(-\frac{1}{\sqrt{2}} \frac{\mu_{s_i}}{\sigma_{s_i}} \right) \right] \quad (21)$$

$$\frac{\partial h}{\partial \sigma_{s_i}^2} = \frac{1}{2} \frac{1}{\sqrt{2\pi} \sigma_{s_i}} \exp \left(-\frac{1}{2} \frac{\mu_{s_i}^2}{\sigma_{s_i}^2} \right) \quad (22)$$

and that

$$\frac{\partial \mu_{s_i}}{\partial \mu_w} = -y_i \begin{bmatrix} \mu_{x_i} \\ 1 \end{bmatrix}, \quad \frac{\partial \mu_{s_i}}{\partial \sigma_w} = \mathbf{0} \quad (23)$$

$$\frac{\partial \sigma_{s_i}^2}{\partial \mu_w} = \begin{bmatrix} 2\sigma_{x_i}^2 \circ \mu_{\tilde{w}} \\ 0 \end{bmatrix}, \quad \frac{\partial \sigma_{s_i}^2}{\partial \sigma_w} = 2\sigma_w \circ \begin{bmatrix} \sigma_{x_i}^2 + \mu_{x_i}^2 \\ 1 \end{bmatrix}. \quad (24)$$

B. Gradient of the Data Generator's Cost

In this section, we turn to the data generator and we focus on approximating $\nabla_{\theta_d} \bar{c}_d$, where \bar{c}_d is defined as in (11). We can separate \bar{c}_d into the sum of n functions acting on each data sample independently, i.e., $\bar{c}_d(\theta_l, \theta_d) = \sum_{i=1}^n \bar{c}_d^{(i)}(\theta_l, \theta_d^{(i)})$, where for each $i \in \{1, \dots, n\}$

$$\bar{c}_d^{(i)}(\theta_l, \theta_d^{(i)}) = \mathbb{E}_{\substack{\mathbf{w} \sim p_l(\cdot; \theta_l) \\ \mathbf{x}_i \sim p_d(\cdot; \theta_d^{(i)})}} \left[\frac{\rho_d}{2} \|\mathbf{x}_i - \hat{\mathbf{x}}_i\|^2 + [1 + y_i(\tilde{\mathbf{w}}^\top \mathbf{x}_i + b)]_+ \right]. \quad (25)$$

By exploiting the properties of the expectation and since $p_d(\cdot; \theta_d^{(i)})$ is an axis-aligned Gaussian distribution with mean μ_{x_i} and standard deviation σ_{x_i} , we can simplify (25) as

$$\bar{c}_d^{(i)}(\theta_l, \theta_d^{(i)}) = \frac{\rho_d}{2} (\|\mu_{x_i} - \hat{\mathbf{x}}_i\|^2 + \|\sigma_{x_i}\|^2) + \mathbb{E}_{\substack{\mathbf{w} \sim p_l(\cdot; \theta_l) \\ \mathbf{x}_i \sim p_d(\cdot; \theta_d^{(i)})}} [[1 + y_i(\tilde{\mathbf{w}}^\top \mathbf{x}_i + b)]_+]. \quad (26)$$

As in the case of the learner, the expectation is a troublesome term having the same form of (14), except for an inverted sign. We adopt the same approximation used in Section III-A to obtain a closed-form function. Accordingly, $t_i = 1 + y_i(\tilde{\mathbf{w}}^\top \mathbf{x}_i + b)$ is assumed to be normally distributed with mean μ_{t_i} and σ_{t_i} . Then, the expectation in (26) can be approximated

as $h(\mu_{t_i}, \sigma_{t_i})$, where function h is defined as in (17). The variance $\sigma_{t_i}^2$ is equal to $\sigma_{s_i}^2$ (16), while μ_{t_i} is given by

$$\begin{aligned} \mu_{t_i} &= \mathbb{E}_{\substack{\mathbf{w} \sim p_l(\cdot; \theta_l) \\ \mathbf{x}_i \sim p_d(\cdot; \theta_d^{(i)})}} [1 + y_i(\tilde{\mathbf{w}}^\top \mathbf{x}_i + b)] \\ &= 1 + y_i(\mu_{\tilde{\mathbf{w}}}^\top \mu_{x_i} + \mu_b). \end{aligned}$$

The samplewise cost of the data generator (26) can thus be approximated as

$$\bar{c}_d^{(i)}(\theta_l, \theta_d^{(i)}) \approx L_d(\mu_{x_i}, \sigma_{x_i}) = \frac{\rho_d}{2} (\|\mu_{x_i} - \hat{\mathbf{x}}_i\|^2 + \|\sigma_{x_i}\|^2) + h(\mu_{t_i}(\theta_d^{(i)}), \sigma_{t_i}(\theta_d^{(i)})). \quad (27)$$

The corresponding gradient is given by

$$\frac{\partial L_d}{\partial \mu_{x_i}} = (\mu_{x_i} - \hat{\mathbf{x}}_i) + \rho_d \left(\frac{\partial h}{\partial \mu_{t_i}} \frac{\partial \mu_{t_i}}{\partial \mu_{x_i}} + \frac{\partial h}{\partial \sigma_{t_i}^2} \frac{\partial \sigma_{t_i}^2}{\partial \mu_{x_i}} \right) \quad (28)$$

$$\frac{\partial L_d}{\partial \sigma_{x_i}} = \sigma_{x_i} + \rho_d \left(\frac{\partial h}{\partial \mu_{t_i}} \frac{\partial \mu_{t_i}}{\partial \sigma_{x_i}} + \frac{\partial h}{\partial \sigma_{t_i}^2} \frac{\partial \sigma_{t_i}^2}{\partial \sigma_{x_i}} \right) \quad (29)$$

where $(\partial h)/(\partial \mu_{t_i})$ and $(\partial h)/(\partial \sigma_{t_i}^2)$ are given as in (21) and (22), and

$$\frac{\partial \mu_{t_i}}{\partial \mu_{x_i}} = y_i \mu_{\tilde{\mathbf{w}}}, \quad \frac{\partial \mu_{t_i}}{\partial \sigma_{x_i}} = \mathbf{0} \quad (30)$$

$$\frac{\partial \sigma_{t_i}^2}{\partial \mu_{x_i}} = 2\sigma_{\tilde{\mathbf{w}}}^2 \circ \mu_{x_i}, \quad \frac{\partial \sigma_{t_i}^2}{\partial \sigma_{x_i}} = 2\sigma_{x_i} \circ (\sigma_{\tilde{\mathbf{w}}}^2 + \mu_{\tilde{\mathbf{w}}}^2). \quad (31)$$

IV. KERNELIZATION

Our game, as in Brückner *et al.* [11], assumes explicit knowledge of the feature space ϕ , where the data generator is assumed to randomize the samples $\hat{\mathbf{x}} = \phi(\mathbf{x})$. However, in many applications, the feature mapping is only implicitly given in terms of a positive semidefinite kernel function $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ that measures the similarity between samples as a scalar product in the corresponding kernel Hilbert space, i.e., there exists $\phi : \mathcal{X} \rightarrow \mathbb{R}$ such that $k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^\top \phi(\mathbf{x}')$. Note that in this setting, the input space \mathcal{X} is not restricted to a vector space like in Section III (e.g., it might contain graphs or other structured entities).

For the representer theorem to hold [39], we assume that the randomized weight vectors of the learner live in the same subspace of the reproducing kernel Hilbert space, i.e., $\mathbf{w} = \sum_j \alpha_j \phi(\hat{\mathbf{x}}_j)$, where $\alpha \in \mathbb{R}^n$. Analogously, we restrict the randomized samples obtained by the data generator to live in the span of the mapped training instances, i.e., $\hat{\mathbf{x}}_i = \sum_{j=1}^n \xi_{ij} \phi(\hat{\mathbf{x}}_j)$, where $\xi_i = (\xi_{i1}, \dots, \xi_{in})^\top \in \mathbb{R}^n$.

Now, instead of randomizing \mathbf{w} and $\hat{\mathbf{X}}$, we let the data generator and the learner randomize $\Xi = (\xi_1, \dots, \xi_n)$ and α , respectively. We assume that the expected costs $\bar{c}_{l/d}$ can be rewritten in terms of α and Ξ in a way that involves only inner products of $\phi(\mathbf{x})$, to take advantage of the kernel trick. This is possible for the term $\mathbf{w}^\top \hat{\mathbf{x}}_i = \alpha^\top \mathbf{K} \xi_i$ in (1) and (2), where \mathbf{K} is the kernel matrix. Hence, the applicability of the kernel trick only depends on the choice of the regularizers. It is easy to see that due to the linearity of the variable shift, existence

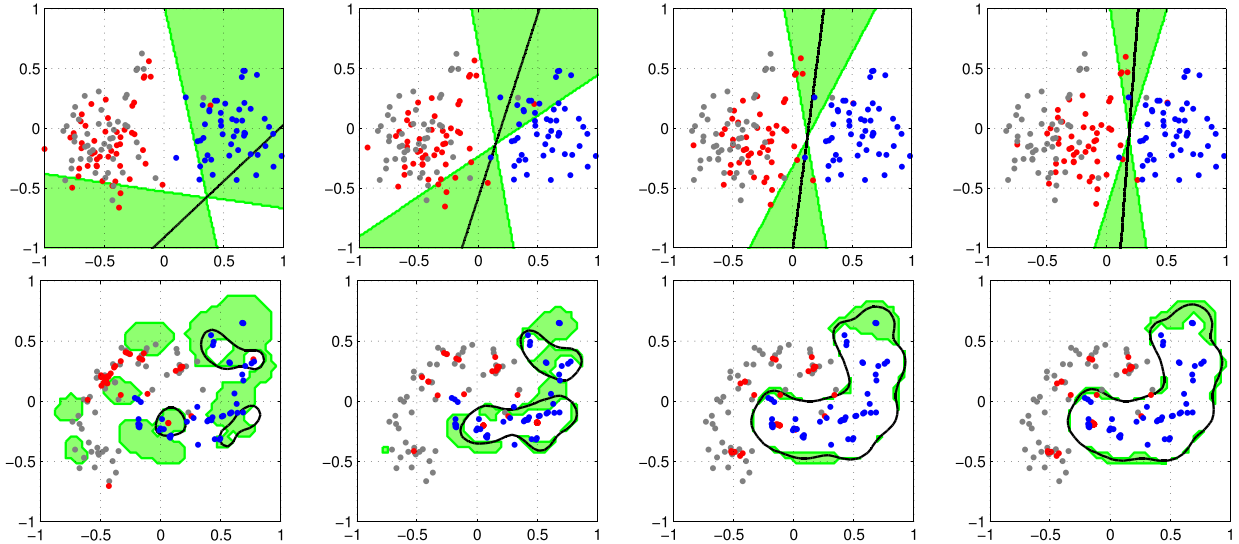


Fig. 1. 2-D examples of randomized prediction games, for SVMs with linear (top) and RBF kernels (bottom). Each row shows how the algorithm gradually converges to a Nash equilibrium. Blue (gray) points represent the legitimate (malicious) class. The mean of each manipulated attack sample is shown as a red point (for clarity, its variance is not shown). The black solid line represents the expected decision boundary, and the green shaded area highlights its variability within one standard deviation. Note how the linear SVM’s decision boundary tends to shift toward the legitimate class, while the nonlinear boundary provides a better enclosing of the same class. This intuitively allows for a higher robustness to different kinds of attack, as it requires the adversary to make a higher number of modifications to the attack samples to evade detection, at the expense of a higher number of misclassified legitimate samples.

and uniqueness of a Nash equilibrium in our kernelized game hold under the same conditions given for the linear case.⁵

Although the data generator is virtually randomizing strategies in some subspace of the reproducing kernel Hilbert space, in reality manipulations should occur in the original input space. Hence, to construct the real attack samples $\{\mathbf{x}_i\}_{i=1}^n$ corresponding to the data generator’s strategy at the Nash equilibrium, one should solve the so-called preimage problem, inverting the implicit feature mapping $\phi^{-1}(\mathbf{K}\xi_i)$ for each sample. This problem is in general neither convex, nor it admits a unique solution. However, reliable solutions can be easily found using well-principled approximations [11], [39]. It is finally worth remarking that solving the preimage problem is not even required from the learner’s perspective, i.e., to train the corresponding, secure classification function.

V. DISCUSSION

In this section, we report a simple case study on a 2-D data set to visually demonstrate the effect of randomized prediction games on SVM-based learners. From a pragmatic perspective, this example also suggests that uniqueness of the Nash equilibrium should not be taken as a strict requirement in our game.

An instance of the proposed randomized prediction game for a linear SVM and for a nonlinear SVM with the RBF kernel is shown in Fig. 1. As one may note from the plots, the main effect of simulating the presence of an attacker that manipulates malicious data to evade detection is to cause the linear decision boundary to gradually shift toward the legitimate class and the nonlinear boundary to find a better enclosing of the legitimate samples. This should generally improve the learner’s robustness to any kind of evasion

attempt, as it requires the attacker to mimic more carefully the feature values of legitimate samples—a task typically harder in several adversarial settings than just obfuscating the content of a malicious sample to make it sufficiently different from the known malicious ones [7], [9].

Based on this observation, any attempt aiming to satisfy the sufficient conditions for uniqueness of the Nash equilibrium will result in an increase of the regularization strength in either the learner’s or the attacker’s cost function. Indeed, to satisfy the condition in Lemma 1, one could sufficiently increase ρ_l , ρ_d , or both. This amounts to increasing the regularization strength of either players, which in turn reduces in some sense their power. Hence, it should be clear that enforcing satisfaction of the sufficient conditions that guarantee the uniqueness of the Nash equilibrium might be counterproductive, by inducing the learner to weakly enclose the legitimate class, either due to a too strong regularization of the learners’ parameters, or by limiting the ability of the attacker to manipulate the malicious samples, thus allowing the learner to keep a loose boundary. This will in general compromise the quality of the adversarial learning procedure. This argument shares similarities with the idea of addressing nonconvex machine learning problems directly, without resorting to convex approximations [37], [38].

Besides improving classifier robustness, finding a better enclosure of the legitimate class may, however, cause a higher number of legitimate samples to be misclassified as malicious. There is indeed a tradeoff between the desired level of robustness and the fraction of misclassified legitimate samples. The benefit of using randomization here is to make the attacker’s strategy less pessimistic than in the case of static prediction games [10], [11], which should allow us to eventually find a better tradeoff between robustness and legitimate misclassifications. This aspect is investigated more systematically in the experiments reported in Section VI.

⁵Note that, on the contrary, manipulating samples directly in the input space would not even guarantee the existence of a Nash equilibrium, as the data generator’s expected cost becomes nonquasi-convex in \mathbf{x} for many (nonlinear) kernels, invalidating Theorem 1.

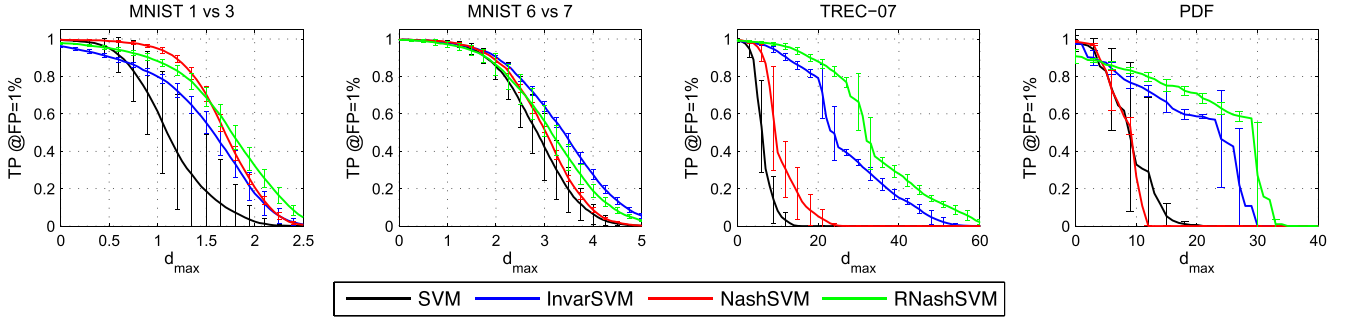


Fig. 2. Security evaluation curves, reporting the average TP at FP = 1% (along with its standard deviation, shown with error bars) against an increasing amount of manipulations to the attack samples (measured by d_{\max}), for handwritten digit (first and second plot), spam (third plot), and PDF (fourth plot) data.

VI. EXPERIMENTS

In this section, we present a set of experiments on handwritten digit recognition, spam filtering, and PDF malware detection. Despite handwritten digit recognition is not a proper adversarial learning task as spam and malware detection, we consider it in our experiments to provide a visual interpretation of how secure learning algorithms are capable of improving robustness to evasion attacks.

We consider only linear classifiers, as they are a typical choice in these settings, and especially in spam filtering [2], [7], [14]. This also allows us to carry out a fair comparison with the state-of-the-art secure learning algorithms, as they yield linear classification functions. We compare our secure linear SVM learner (Section III) with the standard linear SVM implementation [40], and with the state-of-the-art robust classifiers invariant SVM (InvarSVM) [21], [22], and NashSVM [11] (see Section VII).

The goal of these experiments is to test whether these secure algorithms also work well under attack scenarios that differ from those hypothesized during design—a typical setting in security-related tasks, e.g., what happens if game-based classification techniques like that proposed in this paper and NashSVM are used against attackers that exploit a different attack strategy, i.e., attackers that may not act rationally according to the hypothesized objective function? What happens when the attacker does not play at the expected Nash equilibrium? These are important questions to address, as real-world attackers may not play according to the hypothesized objective function.

Security Evaluation: To address the above issues, we consider the security evaluation procedure proposed in [7]. It evaluates the performance of the considered classifiers under attack scenarios of increasing strength. We consider the true positive (TP) rate (i.e., the fraction of detected attacks) evaluated at 1% false positive (FP) rate (i.e., the fraction of misclassified legitimate samples) as performance measure. We evaluate the performance of each classifier in the absence of attack, as in standard performance evaluation techniques, and then start manipulating the malicious test samples to simulate attacks of different strengths. We assume a worst case adversary, i.e., an adversary that has perfect knowledge of the attacked classifier, since we are interested in understanding the worst case performance degradation. Note, however, that other choices are possible, depending on specific assumptions on the adversary’s knowledge and capability [7], [13], [14].

In this setting, we assume that the optimal (worst case) sample manipulation \mathbf{x}^* operated by the attacker is obtained by solving the following optimization problem:

$$\begin{aligned} \mathbf{x}^* \in \arg \min_{\mathbf{x}} & yf(\mathbf{x}; \mathbf{w}) \\ \text{s.t. } & d(\mathbf{x}, \hat{\mathbf{x}}_i) \leq d_{\max} \end{aligned} \quad (32)$$

where y is the malicious class label and $d(\mathbf{x}, \hat{\mathbf{x}}_i)$ measures the distance between the perturbed sample \mathbf{x} and the i^{th} malicious data sample $\hat{\mathbf{x}}_i$ (in this case, we use the ℓ_2 norm, as done by the considered classifiers). The maximum amount of modifications is bounded by d_{\max} , which is a parameter representing the attack strength. It is obvious that the more modifications the adversary is allowed to make on the attack samples, the higher the performance degradation incurred by the classifier is expected to be. Accordingly, the performance of more secure classifiers is expected to degrade more gracefully as the attack strength increases [7], [14].

The solution of the above problem is trivial when we consider linear classifiers, the Euclidean distance, and \mathbf{x} is unconstrained: it amounts to setting $\mathbf{x}^* = \hat{\mathbf{x}}_i - y d_{\max}(\mathbf{w}) / (||\mathbf{w}||)$. If \mathbf{x} lies within some constrained domain, e.g., $[0, 1]$, then one may consider a simple gradient descent with box constraints on \mathbf{x} [9]. If \mathbf{x} takes on binary values, e.g., $\{0, 1\}$, then the attack amounts to switching from 0 to 1 or vice versa the value of a maximum of d_{\max} features which have been assigned the highest absolute weight values by the classifier. In particular, if $y w_k > 0$ ($y w_k < 0$) and the k^{th} feature satisfies $\hat{x}_{ik} = 1$ ($\hat{x}_{ik} = 0$), then $x_k^* = 1$ ($x_k^* = 0$) [7], [14].

Parameter Selection: The considered methods require setting different parameters. From the learners’ perspective, we have to tune the regularization parameter C for the standard linear SVM and InvarSVM, while we have ρ_{-1} and ρ_l for NashSVM and our method, respectively. In addition, the robust classifiers require setting the parameters of their attacker’s objective. For InvarSVM, we have to set K , i.e., the number of modifiable features, while for NashSVM and our method, we have to set the value of the regularization parameter ρ_{+1} and ρ_d , respectively. Furthermore, to guarantee existence of a Nash Equilibrium point, we have to enforce some box constraints on the distribution’s parameters. For the attacker, we restrict the mean of the attack points to lie in $[0, 1]$ (as the considered data sets are normalized in that interval), and their variance in $[10^{-3}, 0.5]$. For the learner, the variance of \mathbf{w} is allowed to vary in $[10^{-6}, 10^{-3}]$, while its mean takes values on $[-W, W]$, where W is optimized together with the

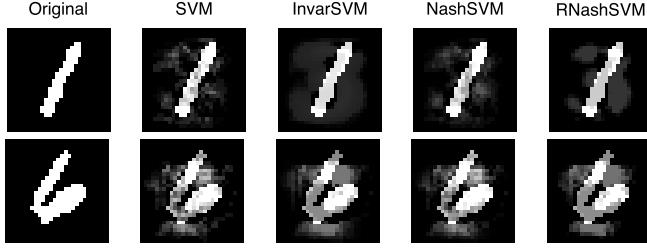


Fig. 3. Examples of obfuscated digits against each classifier when $d_{\max} = 2.5$ for 1 versus 3 (top row), and when $d_{\max} = 5$ for 6 versus 7 (bottom row).

other parameters. All the above-mentioned parameters are set by performing a grid search on the parameter space ($C, \rho_{-1}, \rho_d \in \{0.01, 0.1, 1, 10, 100\}$; $K \in \{8, 13, 25, 30, 47, 52, 63\}$; $\rho_{+1}, \rho_d \in \{0.01, 0.05, 0.1, 1, 10\}$; $W \in \{0.01, 0.05, 0.1, 1\}$), and retaining the parameter values that maximize the area under the security evaluation curve on a validation set. The reason is to find a parameter configuration for each method that attains the best average robustness over all attack intensities (values of d_{\max}), i.e., the best average TP rate at $FP = 1\%$.

A. Handwritten Digit Recognition

Similar to [21], we focus on two two-class subproblems of discriminating between two distinct digits from the MNIST data set [41], i.e., 1 versus 3, and 6 versus 7, where the second digit in each pair represents the attacking class ($y = +1$). The digits are originally represented as grayscale images of 28×28 pixels. They are simply mapped to feature vectors by ordering the pixels in raster scan order. The overall number of features is thus $d = 784$. We normalize each feature (pixel value) in $[0, 1]$, by dividing its value by 255. We build a training and a validation set of 1000 samples each by randomly sampling the original training data available for MNIST. As for the test set, we use the default set provided with this data, which consists of approximately 1000 samples for each digit (class). The results averaged on five repetitions are shown in the first and second plot of Fig. 2, respectively for 1 versus 3, and 6 versus 7. As one may notice, in the absence of attack (i.e., when $d_{\max} = 0$), all classifiers achieved comparable performance (the TP rate is almost 100% for all of them), due to a conservative choice of the operating point ($FP = 1\%$), that should also guarantee a higher robustness against the attack. In the presence of attack, our approach (RNashSVM) exhibits comparable performance with NashSVM on the problem of discriminating 1 versus 3, and to InvarSVM on 6 versus 7. NashSVM outperforms the standard SVM implementation in both cases, but exhibits lower security (robustness) than InvarSVM on 6 versus 7, despite the attacker's regularizer in InvarSVM is not even based on the ℓ_2 norm.

Finally, in Fig. 3, we report two attack samples (a digit from class 1 and one from class 6) and show how they are obfuscated by the attack strategy of (32) to evade detection against each classifier. Notice how the original attacking samples (1 and 6) tend to resemble more the corresponding attacked classes (3 and 7) when natively robust classifiers are used. This visual example confirms the analysis of Section V, i.e., that higher robustness is achieved when the adversary is required to mimic the feature values of samples of the legit-

imate class, instead of slightly modifying the attack samples to differentiate them from the rest of the malicious data.

B. Spam Filtering

In these experiments, we use the benchmark, publicly available, TREC 2007 email corpus [42], which consists of 75 419 real emails (25 220 legitimate and 50 199 spam messages) collected between April and July 2007. We exploit the bag-of-words feature model, in which each binary feature denotes the absence (0) or presence (1) of the corresponding word in a given email [7], [11], [13], [14]. Features (words) are extracted from training emails using the tokenization method of the widely known antispam filter SpamAssassin,⁶ and then $n = 1000$ distinct features are selected using a supervised feature selection approach based on the information gain criterion [43]. We build a training and a validation set of 1000 samples each by randomly sampling the first 5000 emails in chronological order of the original data set, while a test set of about 2000 samples is randomly sampled from the subsequent set of 5000 emails. The results averaged on five repetitions are shown in the third plot of Fig. 2. As in the previous case, in the absence of attack ($d_{\max} = 0$), all the classifiers exhibit a very high (and similar) performance. However, as the attack intensity (d_{\max}) increases, their performance degrades more or less gracefully, i.e., their robustness to the attack is different. Surprisingly, one may notice that only InvarSVM and RNashSVM exhibited an improved level of security. The reason is that these two classifiers are able to find a more uniform set of weights than SVM and NashSVM, and in this case, this essentially requires the adversary to manipulate a higher number of features to significantly decrease the value of the classifier's discriminant function. Note that a similar result has also been heuristically found in [13] and [14].

C. PDF Malware Detection

We consider here another relevant application example in computer security, i.e., the detection of malware in PDF files. The main reason behind the diffusion of malware in PDF files is that they exhibit a very flexible structure that allows embedding several kinds of resources, including Flash, JavaScript, and even executable code. Resources simply consists of keywords that denote their type, and of data streams that contain the actual object, e.g., an embedded resource in a PDF file may be encoded as follows:

```
13 0 obj << /Kids [ 1 0 R 11 0 R ]
/Type /Page... >> end obj
```

where keywords are highlighted in bold face. Recent work has exploited machine learning techniques to discriminate between malicious and legitimate PDF files, based on the analysis of their structure and, in particular, of the embedded keywords [44]–[48]. We exploit here a similar feature representation to that proposed in [45], where each feature denotes the presence of a given keyword in the PDF file. We collected 5993 recent malware samples from the Contagio data set,⁷ and 5951 benign samples from the Web. Following the procedure described in [45], we extracted 114 keywords

⁶<http://spamassassin.apache.org>

⁷<http://contagiodump.blogspot.it>

from the first 1000 samples (in chronological order) to build our feature set. Then, we build training, validation, and test sets as in the spam filtering case, and the average results over five repetitions. Attacks in this case are simulated by allowing the attacker only to increase the feature values of malicious samples, which corresponds to adding the constraint $x \geq \hat{x}_i$ (where the inequality holds for all features) to Problem 32. The reason is that removing objects (and keywords) from malicious PDFs may compromise the intrusive nature of the embedded exploitation code, whereas adding objects can be easily done through the PDF versioning mechanism [9], [46], [48].

The results are shown in the fourth plot of Fig. 2. The considered methods mostly exhibit the same behavior shown in the spam filtering case, besides the fact that, here, there is a clearer tradeoff between the performance in the absence of attack, and robustness under attack. In particular, InvarSVM and RNashSVM are significantly more robust under attack (i.e., when $d_{\max} > 0$) than SVM and NashSVM, at the expense of a slightly worsened detection rate in the absence of attack (i.e., when $d_{\max} = 0$).

To summarize, the reported experiments show that, even if the attacker does not play the expected attack strategy at the Nash equilibrium, most of the proposed state-of-the-art secure classifiers are still able to outperform classical techniques, and, in particular, that the proposed RNashSVM classifier may guarantee an even higher level of robustness. Understanding how this property relates to the use of probability distributions over the set of the classifier's and the attacker's strategies remains an interesting future question.

VII. RELATED WORK

The problem of devising secure classifiers against different kinds of manipulation of samples at test time has been widely investigated in previous work [1], [6], [10], [11], [21]–[27]. Inspired by the seminal work by Dalvi *et al.* [1], several authors have proposed a variety of modifications to existing learning algorithms to improve their security against different kinds of attack. Globerson and Roweis [21] and Teo *et al.* [22] have formulated the so-called InvarSVM in terms of a minimax approach (i.e., a zero-sum game) to deal with worst case feature manipulations at test time, including feature addition, deletion, and rescaling. This paper has been further extended in [23] to allow features to have different *a priori* importance levels, instead of being manipulated equally likely. Notably, more recent research has also considered the development of secure learning algorithms based on zero-sum games for sensor networks, including distributed secure algorithms [27] and algorithms for detecting adversarially corrupted sensors [26].

The rationale behind shifting from zero-sum to nonzero-sum games for adversarial learning is that the classifier and the attacker may not necessarily aim at maximizing antagonistic objective functions. This in turn implies that modeling the problem as a zero-sum game may lead one to design overly-pessimistic classifiers, as pointed out in [11]. Even considering a nonzero-sum Stackelberg game may be too pessimistic, since the attacker (follower) is supposed to move after the classifier (leader), while having full knowledge of the chosen classification function (which again is not realistic in practical settings) [11], [24]. For these reasons,

Brückner *et al.* [10], [11] have formalized adversarial learning as a nonzero-sum game, referred to as static prediction game. Assuming that the players act simultaneously (conversely to Stackelberg games [24]), they devised conditions under which a unique Nash equilibrium for this game exists and developed algorithms for learning the corresponding robust classifiers, including the so-called NashSVM. This paper essentially extends this approach by introducing randomization over the players' strategies.

For completeness, we also mention here that in [25] Bayesian games for adversarial regression tasks have been recently proposed. In such games, uncertainty on the objective function's parameters of either player is modeled by considering a probability distribution over their possible values. To the best of our knowledge, this is the first attempt toward modeling the uncertainty of the attacker and the classifier on the opponent's objective function.

VIII. CONCLUSION

In this paper, we have extended the work in [11] by introducing randomized prediction games. To operate this shift, we have considered parameterized, bounded families of probability distributions defined over the set of pure strategies of either players. The underlying idea, borrowed from [3], [6], and [28], consists of randomizing the classification function to make the attacker select a less effective attack strategy. Our experiments, conducted on an handwritten digit recognition task and on realistic application examples involving spam and malware detection, show that competitive, secure SVM classifiers can be learned using our approach, even when the conditions behind uniqueness of the Nash equilibrium may not hold, i.e., when the attacker may not play according to the objective function hypothesized for her by the classifier. This mainly depends on the particular kind of decision function learned by the learning algorithm under our game setting, which tends to find a better enclosing of the legitimate class. This generally requires the attacker to make more modifications to the malicious samples to evade detection, regardless of the attack strategy chosen. We can thus argue that the proposed methods exhibit robustness properties particularly suited to adversarial learning tasks. Moreover, the fact that the proposed methods may also perform well when the Nash equilibrium is not guaranteed to be unique suggests us that the conditions behind its uniqueness may hold under less restrictive assumptions (e.g., when the SVM admits a unique solution [35], [36]). We thus leave a deeper investigation of this aspect to future work.

Another interesting extension of this paper may be to apply randomized prediction games in the context of unsupervised learning, and, in particular, clustering algorithms. It has been recently shown that injecting a small percentage of well-crafted poisoning attack samples into the input data may significantly subvert the clustering process, compromising the subsequent data analysis [49], [50]. In this respect, we believe that randomized prediction games may help devising secure countermeasures to mitigate the impact of such attacks, e.g., by explicitly modeling the presence of poisoning samples (generated according to a probability distribution chosen by the attacker) during the clustering process.

It is worth finally mentioning that this paper is also slightly related to previous work on security games, in which the goal of the defender is to adopt randomized strategies to protect his or her assets from the attacker, by allocating a limited number of defensive resources, e.g., police officers for airport security and protection mechanisms for network security [51]–[53]. Although our game is not directly concerned to the protection of a given set of assets, we believe that investigating how to bridge the proposed approach within this well-grounded field of study may provide promising research directions for future work, e.g., in the context of network security [52], [53], or for suggesting better user attitudes toward security issues [54]. This may also suggest interesting theoretical advancements, e.g., to establish conditions for the equivalence of Nash and Stackelberg games [51] and to address issues related to the uncertainty on the players' strategies, or on their (sometimes bounded) rationality, e.g., through the use of Bayesian games [25], security strategies, and robust optimization [52], [53]. Another suggestion to overcome the aforementioned issues is to exploit higher level models of the interactions between attackers and defenders in complex, real-world problems, e.g., through the use of replicator equations to model adversarial dynamics in security-related tasks [55]. Exploiting conformal prediction may also be an interesting research direction toward improving current adversarial learning systems [56]. To conclude, we believe that these are all relevant research directions for future work.

REFERENCES

- [1] N. Dalvi, P. Domingos, S. Sanghai, and D. Verma, "Adversarial classification," in *Proc. 10th ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining (KDD)*, Seattle, WA, USA, 2004, pp. 99–108.
- [2] D. Lowd and C. Meek, "Adversarial learning," in *Proc. 11th ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining (KDD)*, Chicago, IL, USA, 2005, pp. 641–647.
- [3] M. Barreno, B. Nelson, R. Sears, A. D. Joseph, and J. D. Tygar, "Can machine learning be secure?" in *Proc. Symp. Inf. Comput. Commun. Secur. (ASIACCS)*, New York, NY, USA, 2006, pp. 16–25.
- [4] A. A. Cárdenas and J. S. Baras, "Evaluation of classifiers: Practical considerations for security applications," in *Proc. AAAI Workshop Eval. Methods Mach. Learn.*, Boston, MA, USA, Jul. 2006, pp. 1–6.
- [5] P. Laskov and M. Kloft, "A framework for quantitative security analysis of machine learning," in *Proc. 2nd ACM Workshop Secur. Artif. Intell. (AISec)*, New York, NY, USA, 2009, pp. 1–4.
- [6] L. Huang, A. D. Joseph, B. Nelson, B. I. P. Rubinstein, and J. D. Tygar, "Adversarial machine learning," in *Proc. 4th ACM Workshop Artif. Intell. Secur. (AISec)*, Chicago, IL, USA, 2011, pp. 43–58.
- [7] B. Biggio, G. Fumera, and F. Roli, "Security evaluation of pattern classifiers under attack," *IEEE Trans. Knowl. Data Eng.*, vol. 26, no. 4, pp. 984–996, Apr. 2014.
- [8] B. Biggio *et al.*, "Security evaluation of support vector machines in adversarial environments," in *Support Vector Machines Applications*, Y. Ma and G. Guo, Eds. Cham, Germany: Springer, 2014, pp. 105–153.
- [9] B. Biggio *et al.*, "Evasion attacks against machine learning at test time," in *Proc. Eur. Conf. Mach. Learn. Principles Pract. Knowl. Discovery Databases (ECML PKDD)*, vol. 8190, 2013, pp. 387–402.
- [10] M. Brückner and T. Scheffer, "Nash equilibria of static prediction games," in *Advances in Neural Information Processing Systems (NIPS)*, vol. 22, Y. Bengio, D. Schuurmans, J. Lafferty, C. K. I. Williams, and A. Culotta, Eds. Cambridge, MA, USA: MIT Press, 2009, pp. 171–179.
- [11] M. Brückner, C. Kanzow, and T. Scheffer, "Static prediction games for adversarial learning problems," *J. Mach. Learn. Res.*, vol. 13, no. 1, pp. 2617–2654, Jan. 2012.
- [12] G. L. Wittel and S. F. Wu, "On attacking statistical spam filters," in *Proc. 1st Conf. Email Anti-Spam (CEAS)*, Mountain View, CA, USA, 2004, pp. 1–16.
- [13] A. Kolcz and C. H. Teo, "Feature weighting for improved classifier robustness," in *Proc. 6th Conf. Email Anti-Spam (CEAS)*, Mountain View, CA, USA, 2009, pp. 1–8.
- [14] B. Biggio, G. Fumera, and F. Roli, "Multiple classifier systems for robust classifier design in adversarial environments," *Int. J. Mach. Learn. Cybern.*, vol. 1, no. 1, pp. 27–41, Dec. 2010.
- [15] M. Christodorescu, S. Jha, S. A. Seshia, D. Song, and R. E. Bryant, "Semantics-aware malware detection," in *Proc. IEEE Symp. Secur. Privacy (SP)*, May 2005, pp. 32–46.
- [16] P. Fogla, M. Sharif, R. Perdisci, O. Kolesnikov, and W. Lee, "Polymorphic blending attacks," in *Proc. 15th Conf. USENIX Secur. Symp. (USENIX-SS)*, Berkeley, CA, USA, 2006, pp. 241–256.
- [17] B. Nelson *et al.*, "Exploiting machine learning to subvert your spam filter," in *Proc. 1st USENIX Workshop Large-Scale Exploits Emergent Threats (LEET)*, Berkeley, CA, USA, 2008, pp. 1–9.
- [18] B. I. P. Rubinstein *et al.*, "ANTIDOTE: Understanding and defending against poisoning of anomaly detectors," in *Proc. 9th ACM Internet Meas. Conf. (IMC)*, New York, NY, USA, 2009, pp. 1–14.
- [19] B. Biggio, B. Nelson, and P. Laskov, "Poisoning attacks against support vector machines," in *Proc. 29th Int. Conf. Mach. Learn.*, 2012, pp. 1807–1814.
- [20] H. Xiao, B. Biggio, G. Brown, G. Fumera, C. Eckert, and F. Roli, "Is feature selection secure against training data poisoning?" in *Proc. 32nd Int. Conf. Mach. Learn. (JMLR W&CP)*, vol. 37, 2015, pp. 1689–1698.
- [21] A. Globerson and S. Roweis, "Nightmare at test time: Robust learning by feature deletion," in *Proc. 23rd Int. Conf. Mach. Learn.*, vol. 148, 2006, pp. 353–360.
- [22] C. H. Teo, A. Globerson, S. Roweis, and A. J. Smola, "Convex learning with invariances," in *Advances in Neural Information Processing Systems (NIPS)*, vol. 20, J. Platt, D. Koller, Y. Singer, and S. Roweis, Eds. Cambridge, MA, USA: MIT Press, 2008, pp. 1489–1496.
- [23] O. Dekel, O. Shamir, and L. Xiao, "Learning to classify with missing and corrupted features," *Mach. Learn.*, vol. 81, no. 2, pp. 149–178, Nov. 2010.
- [24] M. Brückner and T. Scheffer, "Stackelberg games for adversarial prediction problems," in *Proc. 17th ACM Int. Conf. Knowl. Discovery Data Mining (KDD)*, New York, NY, USA, 2011, pp. 547–555.
- [25] M. Großhans, C. Sawade, M. Brückner, and T. Scheffer, "Bayesian games for adversarial regression problems," in *Proc. 30th Int. Conf. Mach. Learn. (JMLR W&CP)*, 2013, vol. 28, no. 3, pp. 55–63.
- [26] K. G. Vamvoudakis, J. P. Hespanha, B. Sinopoli, and Y. Mo, "Detection in adversarial environments," *IEEE Trans. Autom. Control*, vol. 59, no. 12, pp. 3209–3223, Dec. 2014.
- [27] R. Zhang and Q. Zhu, "Secure and resilient distributed machine learning under adversarial environments," in *Proc. 18th Int. Conf. Inf. Fusion*, Jul. 2015, pp. 644–651.
- [28] B. Biggio, G. Fumera, and F. Roli, "Adversarial pattern classification using multiple classifiers and randomisation," in *Proc. 12th Joint IAPR Int. Workshop Struct. Syntactic Statist. Pattern Recognit.*, vol. 5342, 2008, pp. 500–509.
- [29] I. L. Glicksberg, "A further generalization of the Kakutani fixed point theorem, with application to Nash equilibrium points," *Proc. Amer. Math. Soc.*, vol. 3, no. 1, pp. 170–174, Feb. 1952.
- [30] J. B. Rosen, "Existence and uniqueness of equilibrium points for concave N-person games," *Econometrica*, vol. 33, no. 3, pp. 520–534, Jul. 1965.
- [31] D. L. Zhu and P. Marcotte, "Modified descent methods for solving the monotone variational inequality problem," *Oper. Res. Lett.*, vol. 14, no. 2, pp. 111–120, Sep. 1993.
- [32] P. T. Harker and J.-S. Pang, "Finite-dimensional variational inequality and nonlinear complementarity problems: A survey of theory, algorithms and applications," *Math. Program.*, vol. 48, no. 1, pp. 161–220, Mar. 1990.
- [33] C. Geiger and C. Kanzow, *Theorie und Numerik Restringierter Optimierungsaufgaben*. New York, NY, USA: Springer, 1999.
- [34] C. Cortes and V. Vapnik, "Support-vector networks," *Mach. Learn.*, vol. 20, no. 3, pp. 273–297, Sep. 1995.
- [35] S. Abe, "Analysis of support vector machines," in *Proc. 12th IEEE Workshop Neural Netw. Signal Process.*, 2002, pp. 89–98.
- [36] C. J. C. Burges and D. J. Crisp, "Uniqueness of the SVM solution," in *Advances in Neural Information Processing Systems (NIPS)*, S. A. Solla, T. K. Leen, and K.-R. Müller, Eds. Cambridge, MA, USA: MIT Press, 1999, pp. 223–229.
- [37] R. Collobert, F. Sinz, J. Weston, and L. Bottou, "Trading convexity for scalability," in *Proc. 23rd Int. Conf. Mach. Learn. (ICML)*, New York, NY, USA, 2006, pp. 201–208.
- [38] Y. Bengio and Y. LeCun, "Scaling learning algorithms towards AI," in *Large-Scale Kernel Machines*, L. Bottou, O. Chapelle, D. DeCoste, and J. Weston, Eds. Cambridge, MA, USA: MIT Press, 2007.

- [39] B. Schölkopf *et al.*, "Input space versus feature space in kernel-based methods," *IEEE Trans. Neural Netw.*, vol. 10, no. 5, pp. 1000–1017, Sep. 1999.
- [40] C.-C. Chang and C.-J. Lin, "LIBSVM: A library for support vector machines," *ACM Trans. Intell. Syst. Technol.*, vol. 2, no. 3, pp. 1–27, 2011.
- [41] Y. LeCun *et al.*, "Comparison of learning algorithms for handwritten digit recognition," in *Proc. Int. Conf. Artif. Neural Netw.*, 1995, pp. 53–60.
- [42] G. V. Cormack, "Trec 2007 spam track overview," in *TREC*, E. M. Voorhees and L. P. Buckland, Eds. National Institute of Standards and Technology, 2007.
- [43] G. Brown, A. Pocock, M.-J. Zhao, and M. Luján, "Conditional likelihood maximisation: A unifying framework for information theoretic feature selection," *J. Mach. Learn. Res.*, vol. 13, no. 1, pp. 27–66, Jan. 2012.
- [44] C. Smutz and A. Stavrou, "Malicious PDF detection using metadata and structural features," in *Proc. 28th Annu. Comput. Secur. Appl. Conf. (ACSAC)*, New York, NY, USA, 2012, pp. 239–248.
- [45] D. Maiorca, G. Giacinto, and I. Corona, "A pattern recognition system for malicious PDF files detection," in *Machine Learning and Data Mining in Pattern Recognition* (Lecture Notes in Computer Science), vol. 7376, P. Perner, Ed. Heidelberg, Germany: Springer, 2012, pp. 510–524.
- [46] D. Maiorca, I. Corona, and G. Giacinto, "Looking at the bag is not enough to find the bomb: An evasion of structural methods for malicious PDF files detection," in *Proc. 8th ACM Symp. Inf., Comput. Commun. Secur. (ASIACCS)*, New York, NY, USA, 2013, pp. 119–130.
- [47] N. Šrđić and P. Laskov, "Detection of malicious PDF files based on hierarchical document structure," in *Proc. 20th Annu. Netw. Distrib. Syst. Secur. Symp. (NDSS)*, 2013, pp. 1–16.
- [48] N. Šrđić and P. Laskov, "Practical evasion of a learning-based classifier: A case study," in *Proc. IEEE Symp. Secur. Privacy (SP)*, Washington, DC, USA, May 2014, pp. 197–211.
- [49] B. Biggio, I. Pillai, S. R. Bulò, D. Ariu, M. Pelillo, and F. Roli, "Is data clustering in adversarial settings secure?" in *Proc. ACM Workshop Artif. Intell. Secur. (AISec)*, New York, NY, USA, 2013, pp. 87–98.
- [50] B. Biggio *et al.*, "Poisoning complete-linkage hierarchical clustering," in *Structural, Syntactic, and Statistical Pattern Recognition* (Lecture Notes in Computer Science), vol. 8621, P. Fränti, G. Brown, M. Loog, F. Escolano, and M. Pelillo, Eds. Joensuu, Finland: Springer, 2014, pp. 42–52.
- [51] D. Korzhik, Z. Yin, C. Kiekintveld, V. Conitzer, and M. Tambe, "Stackelberg vs. Nash in security games: An extended investigation of interchangeability, equivalence, and uniqueness," *J. Artif. Intell. Res.*, vol. 41, no. 2, pp. 297–327, May 2011.
- [52] T. Alpcan and T. Başar, *Network Security: A Decision and Game-Theoretic Approach*. Cambridge, U.K.: Cambridge Univ. Press, 2010.
- [53] M. Tambe, *Security and Game Theory: Algorithms, Deployed Systems, Lessons Learned*. Cambridge, U.K.: Cambridge Univ. Press, 2011.
- [54] J. Grossklags, N. Christin, and J. Chuang, "Secure or insecure?: A game-theoretic analysis of information security games," in *Proc. 17th Int. Conf. World Wide Web (WWW)*, New York, NY, USA, 2008, pp. 209–218.
- [55] G. Cybenko and C. E. Landwehr, "Security analytics and measurements," *IEEE Security Privacy*, vol. 10, no. 3, pp. 5–8, May 2012.
- [56] H. Wechsler, "Cyberspace security using adversarial learning and conformal prediction," *Intell. Inf. Manage.*, vol. 7, no. 4, pp. 195–222, Jul. 2015.



Samuel Rota Bulò (M'12) received the Ph.D. degree in computer science from the Ca' Foscari University of Venice, Venice, Italy, in 2009.

He has been a Post-Doctoral Researcher with the Ca' Foscari University of Venice since 2013. He is currently a Researcher with the Technologies of Vision Laboratory, Fondazione Bruno Kessler, Trento, Italy. He regularly publishes his research in well-recognized conferences and top-level journals mainly in the areas of computer vision and pattern recognition. He held a research visiting positions

with Instituto Superior Técnico–Technical University of Lisbon, Lisbon, Portugal, the University of Vienna, Vienna, Austria, the Graz University of Technology, Graz, Austria, the University of York, York, U.K., Microsoft Research Ltd., Cambridge, U.K., and the University of Florence, Florence, Italy. His current research interests include areas of computer vision and pattern recognition with particular emphasis on discrete and continuous optimization methods, graph theory, game theory, and field of stochastic modeling.



Battista Biggio (M'07) received the M.Sc. degree (Hons.) in electronics engineering and the Ph.D. degree in electronics engineering and computer science from the University of Cagliari, Cagliari, Italy, in 2006 and 2010, respectively.

He has been with the Department of Electrical and Electronic Engineering, University of Cagliari, since 2007, where he is currently a Post-Doctoral Researcher. In 2011, he visited the University of Tübingen, Tübingen, Germany, where he was involved in the security of machine learning to training data poisoning. His current research interests include secure machine learning, multiple classifier systems, kernel methods, biometrics, and computer security.

Dr. Biggio is a member of the International Association of Pattern Recognition. He serves as a Reviewer for several international conferences and journals.



Ignazio Pillai received the M.Sc. degree (Hons.) in electronics engineering and the Ph.D. degree in electronics engineering and computer science from the University of Cagliari, Cagliari, Italy, in 2002 and 2007, respectively.

He has been with the Department of Electrical and Electronic Engineering, University of Cagliari, since 2003, where he holds a post-doctoral position with the Research Laboratory on Pattern Recognition and Applications. He has authored over 20 publications in international journals and conferences, and

a Reviewer for several international conferences and journals. His current research interests include multilabel classification, multimedia document categorization, and classification with a reject option.



Marcello Pelillo (M'90–SM'04–F'13) is a Professor of Computer Science at Ca' Foscari University of Venice, Italy, where he directs the European Centre for Living Technology (ECLT). He held visiting research positions at Yale University, McGill University, the University of Vienna, York University (U.K.), the University College London, and the National ICT Australia (NICTA). He has published more than 200 technical papers in refereed journals, handbooks, and conference proceedings in the areas of machine learning, computer vision, and pattern

recognition.

Prof. Pelillo is a fellow of the IAPR. He serves (has served) on the Editorial Boards of the journals, such as the IEEE TRANSACTIONS ON PATTERN ANALYSIS AND MACHINE INTELLIGENCE, *IET Computer Vision*, *Pattern Recognition*, and *BRAIN Informatics*, and he serves on the Advisory Board of the *International Journal of Machine Learning and Cybernetics*. He has served (serves) as a Guest Editor for various special issues of the IEEE TRANSACTIONS ON PATTERN ANALYSIS AND MACHINE INTELLIGENCE, *Pattern Recognition*, and *Pattern Recognition Letters*. He has initiated several conferences series as Program Chair (EMMCVPR, IWCW, and SIMBAD) and will serve as a General Chair for ICCV 2017. He has recently been appointed as IEEE SMC Distinguished Lecturer.



Fabio Roli (F'12) received the Ph.D. degree in electronics engineering from the University of Genoa, Genoa, Italy.

He was a Research Group Member of the University of Genoa from 1988 to 1994. He was an Adjunct Professor with the University of Trento, Trento, Italy, from 1993 to 1994. In 1995, he joined the Department of Electrical and Electronic Engineering, University of Cagliari, Cagliari, Italy, where he is currently a Professor of Computer Engineering and Head of the Research Laboratory on Pattern

Recognition and Applications. He was a Very Active Organizer of the international conferences and workshops, and established the popular workshop series on multiple classifier systems. His current research interests include design of pattern recognition systems and their applications.

Dr. Roli is a fellow of the International Association of Pattern Recognition.