

An Adversarial-Risk-Analysis Approach to Counterterrorist Online Surveillance

César Gil Espinasa and Javier Parra-Arnau

Abstract—The Internet, with the rise of the IoT, is one of the most powerful means of propagating a terrorist threat, and at the same time the perfect environment for deploying ubiquitous online surveillance systems. This paper tackles the problem of online surveillance, which we define as the monitoring by a security agency of a set of websites through tracking and classification of profiles that are potentially suspected of carrying out terrorist attacks. We conduct a theoretical analysis in this scenario that investigates the introduction of automatic classification technology compared to the status quo involving manual investigation of the collected profiles. Our analysis starts examining the suitability of game-theoretic-based models for decision-making in the introduction of this technology. We propose an adversarial-risk-analysis (ARA) model as a novel way of approaching the online surveillance problem that has the advantage of discarding the hypothesis of common knowledge. The proposed model allows us to study the rationality conditions of the automatic suspect detection technology, determining under which circumstances it is better than the traditional human-based approach. Our experimental results show the benefits of the proposed model. Compared to standard game theory, our ARA-based model indicates in general greater prudence in the deployment of the automatic technology and exhibits satisfactorily well performance without having to relax crucial hypotheses such as common knowledge and therefore subtracting realism from the problem, although at the expense of higher computational complexity.

Index Terms—adversarial risk analysis, online surveillance, counterterrorism, threat identification, Internet of things.



1 INTRODUCTION

The global threat of terrorism is currently one of the greatest challenges facing our society. Since September 11, Western countries have been allocating more effort and resources to fight terrorism on the national and international scales. However, the resources for the increased security to counter potential large-scale attacks are limited.

In this context, the Internet is one of the most powerful means of propagating a threat with lethal effects, especially in the case of jihadist terrorism. As a matter of fact, a quantitative study [1] of 178 individuals detained in Spain between 2013 and 2016 for activities related to jihadist terrorism shows that there are two crucial factors for understanding their radicalization. On the one hand, face-to-face or online contact with a radicalization agent. On the other hand, the existence of previous social links with other radicalized individuals.

With the rise of the Internet of things (IoT), where billions of online objects embedded in our homes, workplaces and cities will collect and analyze our data, the risk to national security is exacerbated while it opens up a new horizon for more invasive online surveillance technologies.

The revelations by NSA whistleblower Edward Snowden revealed the scale and extent of digital surveillance, particularly on the Internet, by different security and intelligence agencies [2]. In this work, we focus on the problem of *online surveillance* faced by a security agency that monitors

a set of specific websites by tracking and classifying profiles that are potentially suspected of carrying out terrorist attacks. While there is an extensive body of research in decision-making models and risk analysis for fighting terrorism¹, to the best of our knowledge the problem above of online surveillance with counterterrorist purposes, understood as a game between opponents who want to maximize their benefits, has not been tackled yet. Although it is a controversial issue, our interest is to rationalize the matter from a strictly scientific point of view and, in any case, raise new questions and challenges.

The aim of this work is to conduct a theoretical analysis of the rationality conditions implied in the deployment of an online surveillance system for detecting and neutralising potential terrorist threats on the Internet. We consider an approach for evaluating the problem based on adversarial risk analysis (ARA), whose bases are found in [4]. This approach supposes a recently new perspective of decision analysis, providing a robust analytical framework that is a hybrid between game theory and risk analysis. Its objective is to face precisely the risks derived from the intentional actions of intelligent adversaries, which increase security risks, and uncertain results.

We analyse the feasibility of using a technology based on an automatic suspect detection system that covers the functions of investigators who inspect certain websites. That is to say, we aim to determine under which circumstances the tracking and automatic detection model is better than the traditional model (“status quo”) in which the collected user profiles are inspected manually. Our work also allows us to limit the paradox of the false positive [5], which is a controversial problem of mass surveillance systems, since

• César Gil Espinasa is with the Universitat Oberta de Catalunya. Javier Parra-Arnau is with the Department of Computer Science and Mathematics, Universitat Rovira i Virgili, and with the CYBERCAT-Center for Cybersecurity Research of Catalonia, E-08034 Tarragona, Spain.
E-mail: cgile@uoc.edu, javier.parra@urv.cat

our approach is selective and does not infer errors from a broad reference population. Our objective is to carry out a rigorous analysis of the problem.

Next, we summarize the major contributions of this work:

- We analyze the suitability of decision-making models based on standard game theory and ARA, to tackle the problem of online surveillance. Our analysis contemplates the case of sequential defense-attack models, and examines the fulfillment of certain requirements on the defender and attacker’s side.
- We propose an ARA-based model to investigate the problem of online surveillance and analyse the rationality conditions of an automatic threat detection system. Our analysis constitutes a preliminary step towards the systematic application of ARA, in that it aims to establish a point of departure and connection between the analytical framework provided by ARA, a young field within risk analysis, and the problem of online surveillance.
- We conduct an experimental evaluation of the proposed decision-making model and illustrate the typical problem solving approach used in a real case. Our evaluation methodology, in fact, may serve as a template for real problems, which would basically add modelling and computational complexities. Furthermore, we carry out a sensitivity analysis and provide a thorough comparison with a standard game-theoretic approach under assumptions of common knowledge. Our experiments show that our ARA-based model outperforms the standard game-theoretic approach, although at the expense of more costly solutions, from a computational point of view.
- The connection between the ARA models and online counterterrorism sheds new light on the understanding of the suitability of such decision-making models when it comes to applying them to the online surveillance problem. We also hope to illustrate the riveting intersection between the fields of ARA and threat intelligence, in an attempt towards bridging the gap between the respective communities.

The remainder of this paper is organized as follows. Sec. 2 provides some background on online third-party tracking and establishes our assumptions about the surveillance system. Sec. 3 describes the online surveillance problem tackled in this work. Section 4 examines the appropriateness of decision-making models based on standard game theory and ARA, to address the problem of online surveillance. Sec. 5 proposes an ARA-based model for sequential decision-making in the context of online surveillance. Sec. 6 conducts an experimental evaluation of the proposed model. Sec. 7 discusses several aspects of our model in relation to the experimental results. Finally, conclusions are drawn in Sec. 8.

2 BACKGROUND AND ASSUMPTIONS

The purpose of online third-party tracking is behavioral advertising [6], [7], that is to say, showing ads based on the user’s past browsing activity. In this section, we first give a

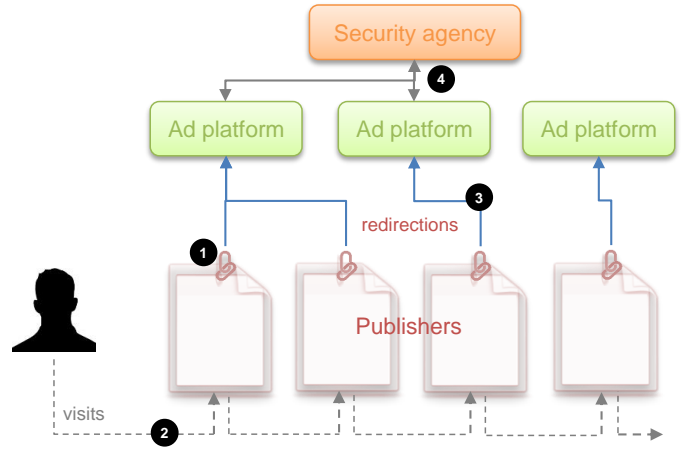


Figure 1: Third-party tracking requires that publishers include a link to the ad platform/s they want to partner with (1). When a user visits pages partnering with this/these ad platform/s (2), the browser is instructed to load the URLs provided by the ad platform/s. Through the use of third-party cookies and other tracking mechanisms, the ad platform/s can track all these visits and build a browsing profile (3). Finally, the information collected by the ad platform/s is shared with the security agency, provided that they have an agreement (4).

brief overview of the main actors of the advertising ecosystem. This will be necessary to understand our assumptions about the online surveillance system assumed in this work, described later in Sec. 2.2.

2.1 Background in Online Third-Party Tracking

The online advertising industry is composed by a considerable number of entities with very specific and complementary roles, whose ultimate aim is to display ads on Web sites. Publishers, advertisers, ad platforms, ad agencies, data brokers, aggregators and optimizers are some of the parties involved in those processes [8]. Despite the enormous complexity and constant evolution of the advertising ecosystem, it is usually characterized in terms of publishers, advertisers and advertising platforms [9], [10], [11], [12], [13].

In this simplified albeit comprehensive terms, the third-party tracking and advertising is carried out as follows. As users navigate the Web and interact with websites, they are observed by both “first parties”, which are the sites the user visits directly, and “third parties”, which are typically hidden trackers such as ad networks embedded on most web pages. The former parties are often known as *publishers* and the latter as *ad platforms*.

Tracking by third-parties begins with publishers embedding in their sites a link to the ad platform/s they want to work with. The upshot is as follows: when a user retrieves one of those Web sites and loads it, their browser is immediately directed to all the embedded links. Then, through the use of third-party cookies, Web fingerprinting or other tracking technologies, the ad platform is able to track the user’s visit to this and any other site partnering with it. Third parties can learn not only the Web pages visited and hence its content, but also the user’s location through their IP address, and, more importantly, their Web-browsing interests, also known as *navigation trace*.

2.2 Assumptions

In this section, we describe our assumptions about the surveillance system deployed by a security agency for detecting possible terrorist threats on websites of interest. The practical details of this system and possible anti-tracking countermeasures, however, are beyond the scope of this work. The purpose of our analysis is not to go deeply into these details but rather to study the rationality conditions of deploying such an online surveillance technology.

First, we suppose that a security agency wants to develop a Web infrastructure on which to apply an online automatic threat detection system. The websites or publishers targeted by the agency will be those that make it possible to detect threats. For example, certain web forums where ISIS recruiting messages appear with certain frequency are sites that are susceptible to being investigated by the security agency.

In addition, we suppose that it is possible to track users' activity in the target sites, or in other words, there are advertising and tracking companies operating in these sites. We acknowledge, however, that there may be sites such as those hosted on the Dark web or others that are on the Internet that cannot be subject to surveillance because there are no ad platforms and tracking companies.

We assume that the agency can contract the services of the trackers available at the target sites to capture the users' visit data, which may include, among others, their activity within the site, location, IP address and Web-browser fingerprints. Once properly treated, all such data may allow the agency to reidentify a given Web user, possibly with the help of the Internet service provider in question.

In essence, the infrastructure assumed is based on three well-differentiated activities. In a first stage, the agency selects its target publishers and hires the services of the companies that track them to obtain the users' raw visit data. In a second optional stage, the agency exploits the data captured by the third-party trackers through an automatic system based on artificial intelligence methods (classifiers) so that, once the navigation trace of each user is extracted, it is possible to obtain a binary classification: suspicious or not suspicious. The threat detection algorithm that underlies this automatic system inevitably has certain sensitivity and specificity parameters (false positives). In a third and final stage, whether the automatic system has been deployed or not, there is an essential manual investigation of the flagged users by security experts. It should be noted that this type of architecture has two types of limited resources that are well differentiated: resources for hiring trackers and resources for the manual investigation of the collected profiles. In this work, we consider the cost of first type of resources is negligible compared to that of the latter. Fig. 1 provides a conceptual depiction of the surveillance infrastructure assumed in this work.

3 THE PROBLEM OF ONLINE SURVEILLANCE

In this section, we describe the problem of online surveillance from the intrusion-detection problem posed and

Table 1: Parameters of the online surveillance problem.

Symbol	Description
α	Probability of ASC alarm due to suspicion (true positive)
β	Probability of ASC alarm without suspicion (false positive)
π	Probability of presence of a suspicious user
ρ	Probability of manual investigation without using ASC
ρ_1	Probability of manual investigation when the ASC generates an alarm
ρ_0	Probability of manual investigation when the ASC does not generate an alarm
c	Cost of manual investigation; $c \leq \phi d$, $\phi \leq 1$
d	Damage derived from an undetected suspect
ϕ	Cost/damage coefficient of the system
b	Benefit for suspects not detected; $l \geq (1 + \lambda)b$, $\lambda \leq 1$
l	Loss for suspects not detected
λ	Benefit/loss coefficient of the suspect

solved by [14]². It is also appropriate to point out the work of Merrick and McLalay [16] on the use of scanners against smuggling of nuclear devices in cargo containers. Both works treat physical or logical security problems and assess their conditioning factors under uncertainty with the use of automatic threat detection systems. We rely on the cited works to define the problem at hand.

Suppose we are going to give support to the decision making of a security agency that has jurisdiction in a territory to prevent and neutralize attacks perpetrated by terrorists who use the Internet as a resource for carrying out their attacks. In general terms, we assume that the agency wishes to suffer the least possible harm, and, on the contrary, the terrorists want to cause the greatest possible damage. Faced with a normal Internet user, we define the suspect as a user whose digital activity can be considered a threat that must be investigated by the agency.

Suppose then that, in a certain period of time, the security agency will carry out online surveillance tasks on a series of websites that, based on expert knowledge, have been classified as susceptible to being used (propaganda, training, forums, etc.) by users who could potentially acquire the capabilities to prepare and/or carry out attacks.

To monitor these sites, the security agency uses a digital technology based on automatic detection of threats, described in Sec. 2, which consists of two well-defined complementary functions: automatic collection and classification of user profiles. Firstly, and based on tracking the digital activity of users who browse the target websites, the system collects the navigation traces, which result in unique user profiles. Secondly, of the profiles collected, the system is able to detect those that are potentially suspicious with certain sensitivity and specificity rates. More specifically, these classifiers are based on artificial intelligence methods. The use of the classifier is optional and in any case the system can always be supported by an "ad hoc" manual

2. According to [15], intrusion detection systems (IDSs) are hardware or software systems that automate the process of monitoring events that have occurred in a computer system or network, analysing them to detect security problems.

investigation by experts whose criteria we will assume to be totally reliable. The classifier analyses each profile and if it considers it to be suspicious, it generates an alarm signal. Afterwards, the agency decides whether or not to investigate the profile based on available (limited) resources. Therefore, the agency makes decisions about whether or not to investigate according to the state (signal or lack of signal) of the system. However, when the system generates a signal, the agency does not know with certainty whether it is a real threat or whether the system has generated a false alarm. On the other hand, the suspect user's main objective is not to be detected by the surveillance system, which would imply, immediately and to simplify, the success of their actions.

The aim of the agency is to configure the system by choosing a point in its effectiveness function that minimizes the total cost of surveillance (the cost is not necessarily a monetary value but we can treat values such as image, privacy, etc., or in any case monetize them). Thus, we initially define the probability of detection α as the probability of classifying a suspect conditioned on the user really being a suspect, and the probability of a false positive β as the probability of classifying a suspect conditioned on the user not being a suspect. In a perfect surveillance system, we would suppose $\alpha = 1$ and $\beta = 0$. However, and in general, online surveillance technology is such that a high value of α also implies a high value of β , due to the variability of the data associated with the normal and abnormal traces and the imprecision of the algorithms used by these types of systems.

In general terms, the navigation trace of potential suspects will depend on factors such as the benefit derived from acquiring the capacities to carry out terrorist acts of different levels; the loss that they will receive if they are captured; and the probability that they will be detected. We assume that a potential terrorist obtains a benefit b if their navigation is not detected. If it is detected, the user incurs a loss l over a non-positive net benefit of $(b - l) \leq 0$. Suppose that it is reasonable to think that $l = (1 + \lambda)b$, with $\lambda \leq 0$. The loss can take different forms depending on the nature of the terrorist potential (cost of legal persecution, reputation, intimidating effect, etc.). We denote by π the probability of the presence of a suspicious user in the set of monitored sites.

The agency complements the system with a manual investigation conducted by security experts. In general, it is expensive to always carry out manual investigations (it is obvious that it is a limited resource). When the agency does not deploy the automatic system, expert investigators must manually investigate a proportion ρ of the user profiles collected. When the system is deployed, experts can only investigate a proportion ρ_1 of the profiles that generated alarm signals and a proportion ρ_0 of the profiles that did not generate signals. The agency incurs a cost c every time the experts conduct a manual investigation. We assume that expert manual investigation always confirms or discards threats with certainty (it is 100% effective). If the agency detects a threat it will not incur any loss other than the cost c of the manual investigation. When a suspicious profile is not detected, the agency incurs a damage d . Suppose again that it is reasonable to think that $c \leq \phi d$, with $\phi \leq 1$. It is usual to estimate these possible damages in the risk

assessment phase before implementing and configuring the detection system. Traditionally, the quality function of a detection system is modeled through its relative operating characteristic (ROC) curve, although other evaluation functions can be appropriate as we will see in the next section. Table 1 summarises the parameters of the problem of online surveillance defined in this section.

4 ANALYSIS OF DECISION-MAKING MODELS

The terrorist attacks occurred in Western countries in the last decades have sparked a growing interest in decision-making models and risk analysis for fighting terrorism. We refer the reader to [3] for a complete review of the field.

The vast majority of this literature adopts a game-theoretic approach [17]. Examples comprise [18], which studies multiattribute utility functions for the defender and attacker, and for simultaneous and sequential actions, to compute Nash equilibria; and [19] which proposes several max-min optimization models to tackle defender-attacker, attacker-defender and defender-attacker-defender problems. A hybrid model between game theory and risk analysis is ARA [4], a recently new perspective of decision analysis that differs from standard game theory in that it makes no assumptions of common knowledge.

The other mainstream literature adopts a decision-analysis approach. Among such works, we highlight [20], which uses decision trees to assess man-portable air defense systems countermeasures. The recurrent problem of decision analysis, however, is the need to evaluate the likelihood of the actions of the others, which is a central issue of the Bayesian approach to games³.

In this section, we focus on standard game theory and ARA, and analyze their suitability to tackle the online surveillance problem described in Sec. 3. Since the aim of our analysis is to gain insight into the rationality of online surveillance with the principle of being as close as possible to reality, we define the following requirements for such a model:

- both opponents (intelligent, rational) want to maximize their utility;
- there is uncertainty about the attacker's actions due to uncertainty about their utilities and probabilities;
- the information on the evaluation of the objectives between opponents is incomplete, with the possibility of obtaining it partially through different sources that we will call intelligence (experts, historical data and/or statistical distributions);
- and it is possible to model simultaneous and non-simultaneous (sequential) decisions.

Throughout this section, we shall follow the convention of using uppercase letters for random variables (r.v.'s), and lowercase letters for the particular values they take on. Accordingly, \hat{p} will denote approximation, estimation, as a result of Monte Carlo simulation; and $p^k \sim P$ will denote the former is the k -th iteration of the Monte Carlo simulation of the latter r.v. In text, we shall drop the superindex k for notational simplicity.

3. We would like to stress that the tension between game-theoretic and decision-analytic approaches to decision-making problems with adversaries is by no means exclusive of counterterrorism models [17].

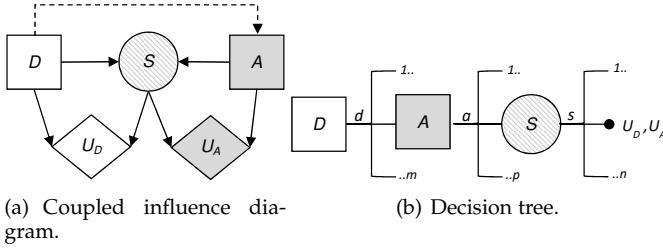


Figure 2: Sequential defence-attack model.

4.1 Sequential Defence-Attack Model

To study the appropriateness of standard game theory and ARA, we develop first the sequential defence-attack model, which is one of the two standard counterterrorism model formulations⁴. We will use this model to analyze the problem that is the objective of this work. For the sake of comparison, we consider the following example of counterterrorism scenario.

Example 1 (Counterterrorism scenario). The authority of an airport (D , the defender) must decide whether or not to install body scanners at the security checkpoints of an airport, replacing the X-ray scanners. On the other hand, a terrorist group (A , the attacker) decides whether or not to try to smuggle a bomb onto an airplane. D makes the first move, so A can see if the new body scanners are in use when they arrive at the airport. Because A can observe the actions of D before deciding their move, they do not need to know their probabilities or utilities. But D must have a distribution for A , which specifies its utilities and probabilities.

In this model, the defender makes the first move, deploys their defensive resources and makes a certain choice in order to position themselves against the possible terrorist threat. The attacker, after having observed this decision, evaluates their options and carries out an attack.

We assume that the defender initially has a discrete set of possible decisions $D = \{d_1, d_2, \dots, d_m\}$ and that the attacker can respond with one of their possible attacks $A = \{a_1, a_2, \dots, a_p\}$. As a consequence of these actions, a result is produced. This result is the only uncertainty of the problem and depends probabilistically on $(d, a) \in D \times A$. The decision sets can include the option to do nothing or combine several defences or several attacks. To simplify the discussion, we consider only two possible values for the result, $S = \{0, 1\}$, which represents the failure or success of the attack. Thus, the defender and the attacker can have different probability distributions for the possibility of success, given a pair (d_i, a_j) . They can also have different utility functions.

To visualize the situation, we have built the influence diagram and the decision tree corresponding to the problem at hand. These are two decision analysis tools that help us to gain a clearer view of the sequential decisions that have to be made.

An influence diagram is a directed acyclic graph with three kinds of nodes: decision nodes, which are shown as squares; chance or uncertainty nodes, shown as circles; and value nodes, shown as hexagons. In addition, an influence

diagram can have three types of arcs depending on their destination: if the arc arrives at a decision node, this indicates that the decision is made knowing the value of the predecessor; if it arrives at a random node, then the uncertainty depends on the predecessor node conditioned probability; and if it arrives at a value node, then the utility reflected in that value node depends on the values of the predecessors.

Fig. 2 shows the coupled influence diagram of the model. In view of this, we assume that the consequences for the defender and the attacker depend, respectively, on (d, s) and (a, s) . This is then inferred to the defender's utility node by two arcs coming from the decision node, D , and the uncertainty node, S , which represents the result. The decision node representing the attacker's utility also has two arcs, which in this case come from the decision node A and the uncertainty node S . It is also reflected that the result node, S , depends, in this case probabilistically, on the defender's initial action and the attacker's response. The influence diagram arc from the node of the defender's first decision to the attacker's node reflects that the defender's choice is observed by the attacker before they decide on their attack.

We also show in Fig. 2(b) the decision tree for this problem, clearly reflecting its sequential nature. First, a decision is made corresponding to the set D ; once the attacker observes this decision, they decide whether to attack or not; the final result is produced as a consequence of these two actions. Note that there are two utility values in the terminal node of the tree. Each of these represents the utility that corresponds to each of the actors: one value refers to the defender's utility and the other value refers to the attacker's utility. The fact that there are several branches of each of the nodes refers to the possible decisions or results, which is the case of the chance node, which can be taken in each of them. The number of possible decisions in each decision node is not always the same and that is reflected in the decision tree.

4.2 Analysis based on Standard Game Theory

The focus of game theory to solve the posed problem requires obtaining the utility functions of the defender $u_D(d, s)$ and attacker $u_A(a, s)$, as well as evaluating the probability of the event $S|d, a$ for each of the participants, which we designate as $p_D(S|d, a)$ and $p_A(S|d, a)$ for the defender and the attacker respectively. Standard game theory requires as initial assumption that the defender knows the attacker's utilities and probabilities and the attacker knows the defender's utilities and probabilities, this being *common knowledge*. If this happens, a solution to the problem can be obtained from the decision tree (Fig. 2(b)) by backward induction as follows.

In node S , it is common knowledge for the two participants both the defender's expected utility associated with each pair $(d, a) \in D \times A$,

$$\begin{aligned} \psi_D(d, a) = & p_D(S = 0|d, a) u_D(d, S = 0) \\ & + p_D(S = 1|d, a) u_D(d, S = 1), \end{aligned}$$

and the attacker's expected utility associated with $(d, a) \in D \times A$,

4. The other is the simultaneous defend-attack model.

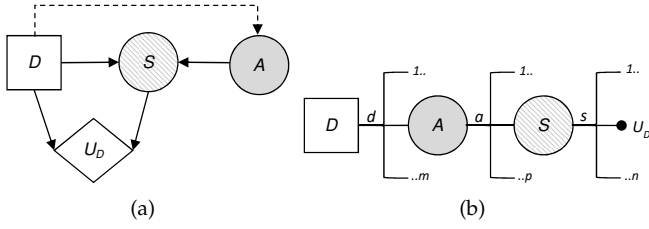


Figure 3: (a) Influence diagram of the defender and (b) decision tree of the defender.

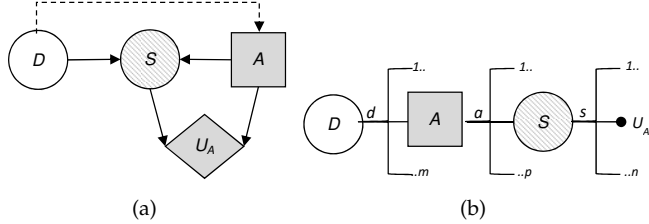


Figure 4: (a) Influence diagram of the attacker and (b) decision tree of the attacker.

$$\begin{aligned} \psi_A(d, a) = & p_A(S = 0 | d, a) u_A(d, S = 0) \\ & + p_A(S = 1 | d, a) u_A(d, S = 1). \end{aligned}$$

Knowing what the defender will do in decision node D , the attacker can determine what is their best attack in node A , after observing the defensive action of the defender, solving the problem

$$a^*(d) = \arg \max_{a \in A} \psi_A(d, a), \forall d \in D.$$

This is also known by the defender due to the hypothesis of common knowledge. The defender can determine their best decision in the decision node D , solving the problem

$$d^* = \arg \max_{d \in D} \psi_D(d, a^*(d)).$$

Thus, under the assumption of common knowledge, standard game theory predicts that the defender will choose $d^* \in D$ in node D ; then the attacker will respond by choosing the attack $a^*(d^*)$. The pair $(d^*, a^*(d^*))$ determines a solution of the game and is a Nash equilibrium.

4.3 Analysis based on ARA

Now we abandon the assumption of common knowledge. It should be taken into account that ARA serves here as support to the defender. To do this, we treat the attacker's behaviour in node A as uncertainty from the defender's point of view and we model this added uncertainty. This is reflected in the influence diagram and the decision tree, as the attacker's decision node has become a chance node, replacing the square by a circle. Looking at the influence diagram (Fig. 3(a)) we need now to obtain $p_D(A|d)$, the probability that the defender will assign to the attack what the attacker will choose once they have observed every defensive move $d \in D$ of the defender. The defender also needs to evaluate $u_D(d, s)$ and $p_D(S|d, a)$, already described above.

Once these data have been evaluated, the defender can solve their decision problem with backward induction considering the decision tree (Fig. 3(b)). Then, the defender will obtain their expected utility in the node S , $\psi_D(d, a)$, for

each pair $(d, a) \in D \times A$ in the same way as in the previous approach. It is at this time when the defender can use the evaluation of the probability of what the attacker will do faced with each of the defender's decisions, $p_D(A|d)$, to determine their expected utility in the node A for each $d \in D$, with the expression

$$\psi_D(d) = \sum_{i=1}^p \psi_D(d, a_i) \hat{p}_D(A = a_i | d).$$

Finally, the defender can find the decision that maximizes their expected utility in node D , solving the problem

$$d^* = \arg \max_{d \in D} \psi_D(d, a^*(d)).$$

Therefore, the best strategy for the defender for the defence-attack model is to choose first d^* in node D after having observed $s \in S$.

The key now is how to evaluate $p_D(A = a | d)$. To do this, ARA assumes the defender can use a statistical method if they have historical data on the attacker's behaviour in similar situations. To complement this evaluation, the defender could also use expert opinions. However, as we shall describe in Sec. 5, an approach could be modelling the uncertainty that the defender has about the attacker's decision assuming (i) that the attacker wants to maximize their expected utility; and (ii) that that the defender's uncertainty in evaluating this probability stems from the uncertainty that the defender has about the attacker's probabilities and utilities associated with the attacker's decision problem. In short, the evaluation is limited to analyzing the attacker's decision problem from the defender's point of view (Figures 5 and 6). The evaluation of the attacker's probabilities and utilities from the defender's perspective will be based on all the information that the defender has available, which can include previous data from similar situations and expert opinions. If the defender does not have this kind of information, they can use an uninformative or reference distribution to describe $p_D(A = a | d)$. Therefore, to obtain $p_D(A = a | d)$, the defender needs to evaluate $u_A(a, s)$ and $p_A(S|d, a)$, the attacker's utilities and probabilities, which are unknown to the defender.

If the defender can access the attacker's probabilities and utilities they will learn, by the same procedure as in the game theory approach, the action that the attacker would give most expected utility, $a^*(d)$, for each $d \in D$, and therefore, $p_D(A = a^*(d) | d) = 1$. This would imply that the attacker's decision would be anticipated by the defender, and therefore they would not need to evaluate $p_D(A = a | d)$.

We start, therefore, from the assumption that the defender does not know these two quantities, but can recognize their uncertainty about them by means of a probability distribution $F = (U_A(a, s), P_A(S|d, a))$ and solve the attacker's decision problem using backward induction on the decision tree of Fig. 4(b) with the expression

$$\begin{aligned} \Psi_A(d, a) = & P_A(S = 0 | d, a) U_A(a, S = 0) \\ & + P_A(S = 1 | d, a) U_A(a, S = 1). \end{aligned}$$

In node A , assuming that the attacker wants to maximize their expected utility, the defender's distribution on the

Requirements	Standard game theory	ARA
opponents aim to maximize their utility	✓	✓
uncertainty about the attacker's actions	✗	✓
incomplete information about the evaluation of the objectives between opponents	✗	✓
simultaneous and sequential decisions	✓	✓

Table 2: Comparison between standard game theory and ARA.

attacker's choice when the defender has considered their defense d is

$$p_D(A = a^* | d) = P_F[a^* = \arg \max_{a \in A} \Psi_A(d, a)].$$

This distribution can be approximated using Monte Carlo simulation methods, generating n values, such that

$$\hat{p}_D(A = a | d) = |\{a = \arg \max_{x \in A} \psi_A^i(d, x)\}|/n, \forall a \in A.$$

Once the defender has completed their evaluations, from these data they can solve their problem in the S node for each $(d, a) \in D \times A$ with the expression

$$\Psi_D(d, a) = p_D(S = 0 | d, a) u_D(d, S = 0) + p_D(S = 1 | d, a) u_D(d, S = 1).$$

Then their estimated expected utility is

$$\hat{\psi}_D(d) = \sum_{i=1}^k \psi_D(d, a_i) \hat{p}_D(A = a_i | d)$$

and finally their optimal decision is

$$d^* = \arg \max_{d \in D} \hat{\psi}_D(d).$$

In view of our analysis of standard game theory and ARA, we regard the latter as the most appropriate model for evaluating the online surveillance problem defined in Sec. 3. The counterterrorism modelling based on common-knowledge assumptions entails that parts have too much information about their counterparts, which does not seem to make sense in a field in which secrecy tends to be an advantage. In a scenario where the adversary wishes to increase the risks of the defender, it seems unusual that the defender will have a full knowledge of their objectives, intentions or possible movements. Similarly, it seems unrealistic that the adversary fully knows the objectives, intentions or possible movements of the defender. Table 2 summarizes the analysis of the two models by matching the initial requirements defined at the beginning of Sec. 4.

5 AN ARA MODEL FOR THE ONLINE SURVEILLANCE PROBLEM

We present an ARA model to evaluate the problem of online surveillance described in Sec. 3. This model will allow us to analyse the rationality conditions of the automatic threat detection system.

We assume that we support an agent (the agency, the defender, D) in their decision-making in relation to deploying an online surveillance system to monitor a set of selected

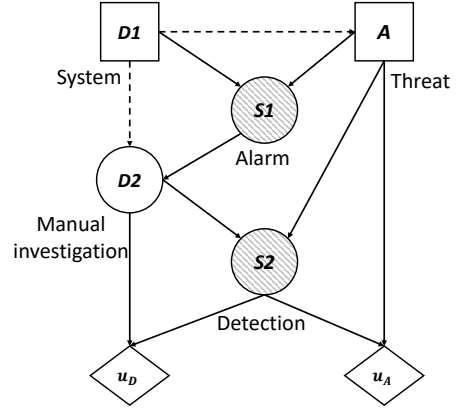


Figure 5: Influence diagram for the online surveillance problem.

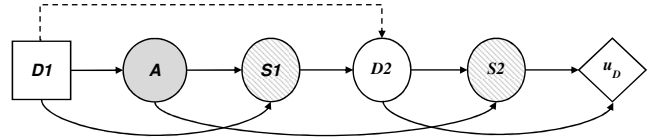


Figure 6: Influence diagram for the defender's decision problem.

websites, faced with the threat posed by the presence of the other agent (the suspect, the adversary, A) in the target sites. We assume that both agents operate monolithically.

According to the premises described in Sec. 2, we assume that the dynamics of the defender and the adversary can be described by means of a sequential defence-attack decision model represented in Fig. 5 as an influence diagram coupled for the two agents.

To begin and given a set of target websites, the defender makes his initial decision $d_1 = \{0, 1\}$ (0 is No, 1 is Yes) about using the technology, represented by the decision node D_1 . The adversary knows about these tracking measures and even so decides to be present in the set of sites, $a = \{0, 1\}$, represented by the decision node A . The automatic system, in the case that it is deployed ($d_1 = 1$), can lead to a system alarm signal, $s_1 = \{0, 1\}$, represented by the node of uncertainty S_1 shared by the defender and the adversary (if the system is not deployed, $s_1 = 0$ unfaillingly). Depending on the previous result, the defender manually investigates the alarm, $d_2 = \{0, 1\}$, represented by the node of uncertainty D_2 , to the degree that their (limited) resources allow. All this leads to the final result of the success/failure of the two agents, $s_2 = \{0, 1\}$, represented by the node of uncertainty S_2 . We understand as success for the security agency the fact of detecting the threat and avoiding its potential actions, and failure is understood as the opposite. For the adversary, success and failure are the reverse events of the defender.

The utility u_D obtained by the defender depends on the added cost of the manual investigation and the final success of the surveillance (nodes D_2 and S_2), on which their utility function is applied. Similarly, the utility u_A obtained by the attacker depends on the added cost of access to the set of sites and the final success of the surveillance (nodes A y S_2), on which their utility function is applied.

5.1 The Defender's Decision Problem

We describe in this section the defender's decision problem, illustrated by an influence diagram in Fig. 6, where the

threat appears as a probability node A , from the point of view of D , which, given a collected profile should:

- decide if they use the technology, assigning values $d_1 = \{0, 1\}$ in node D_1 ;
- face the possible existence of a threat $a = \{0, 1\}$ in node A ;
- observe optionally, given the case, the result of the automatic detection system, $s_1 = \{0, 1\}$, in node S_1 .
- establish proportions of profiles investigated manually based on the available resources, assigning values $d_2 = \{0, 1\}$ in node D_2 ;
- observe the final result of the surveillance, $s_2 = \{0, 1\}$, in node S_2 ;
- add their costs and evaluate the results with their utility function u_D .

To solve the decision problem, D , it is necessary to evaluate the probability distributions, $p_D(A|d_1)$, $p_D(S_1|a, d_1)$, $p_D(D_2|d_1, s_1)$ and $p_D(S_2|a, d_2)$ and the utility function $u_D(d_2, s_2)$. Assuming that D is capable of providing such inputs, we would proceed by applying standard decision analysis calculations based on dynamic programming to obtain the optimal decision.

- First, for each relevant scenario (d_2, s_2) add the consequences and obtain the utility $u_D(d_2, s_2)$.
- In node S_2 , calculate the expected utilities:

$$\psi_D(d_1, s_1, d_2) = \sum_{s_2} u_D(d_2, s_2) p_D(S_2 | d_2, a).$$

- In node D_2 , calculate the expected utilities:

$$\psi_D(d_1, s_1) = \sum_{d_2} \psi_D(d_1, s_1, d_2) p_D(D_2 | d_1, s_1).$$

- In node S_1 , calculate the expected utilities:

$$\psi_D(d_1, a) = \sum_{s_1} \psi_D(d_1, s_1) p_D(S_1 | d_1, a).$$

- In node A , calculate the expected utilities:

$$\psi_D(d_1) = \sum_{d_2} \psi_D(d_1, a) p_D(A | d_1).$$

- Finally, the decision node D_1 maximizes the expected utility and stores the corresponding optimal initial decision d_1^* .

$$\psi_D = \arg \max_{d_1} \psi_D(d_1).$$

Then, d_1^* describes the optimal decision for the defender.

It should be kept in mind that we can describe the defender's optimization problem with the expression

$$d_1^* = \arg \max_{d_1} \sum_a \sum_{s_1} \sum_{d_2} \sum_{s_2} u_D(d_2, s_2) p_D(S_2 | d_2, a) \times p_D(D_2 | d_1, s_1) p_D(S_1 | d_1, a) p_D(A | d_1).$$

Note that of the four values required by the agency, $p_D(A|d_1)$ is the most problematic, insofar as it involves the defender's beliefs about the adversary's decision once they have observed the defender's initial decision d_1 . This is an evaluation that requires strategic thinking for which we propose an approach based on ARA. For this we need to solve the adversary's decision problem, assuming uncertainty about their evaluations and propagating it to obtain

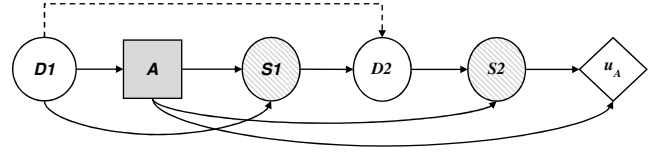


Figure 7: Influence diagram for the adversary's decision problem.

its expected distribution based on the optimal presence of the adversary in the set of monitored sites. We discuss this in the following section.

5.2 The Adversary's Decision Problem

We describe the dynamics of the threat, illustrated as an influence diagram in Fig. 7, according to the defender's point of view, where D_1 is now a probability node for the attacker, who must:

- observe the initial decision of D , $d_1 = \{0, 1\}$;
- decide on their presence in the set of monitored sites, $a = \{0, 1\}$, with impact over time if they are not detected;
- observe their success $s_2 = \{0, 1\}$ after the defender makes their allocations $d_2 = \{0, 1\}$ on the manual investigation of the profiles;
- and finally, add their costs and obtain the corresponding utility u_A .

In order to solve the decision problem, we assume that the adversary wants to maximize their expected utility. They therefore need to evaluate $p_A(S_1|a, d_1)$, $p_A(D_2|s_1, d_1)$, $p_A(S_2|d_2, a)$, and $u_A(a, s_2)$. We will not easily have these values so we model the defender's uncertainty about them with random probability distributions. Then we can propagate this uncertainty using the standard reduction algorithm of influence diagrams and obtain the optimal and random decision $a = \{0, 1\}$ for each value of $d_1 = \{0, 1\}$. This will provide us with the required distribution $p_D(A | d_1) = P(A^*(d_1) = a)$.

- Add the consequences and obtain the random utility $U_A(a, s_2)$, for each (a, s_2) .
- In node S_2 , calculate the expected random utilities:

$$\Psi_A(a, s_1, d_2) = \sum_{s_2} U_A(a, s_2) P_A(S_2 | d_2, a).$$

- In node D_2 , calculate the expected random utilities:

$$\Psi_A(d_1, a, s_1) = \sum_{d_2} \Psi_A(a, s_1, d_2) P_A(D_2 | d_1, s_1).$$

- In node S_1 , calculate the expected random utilities:

$$\Psi_A(d_1, a) = \sum_{s_1} \Psi_A(d_1, a, s_1) P_A(S_1 | d_1, a).$$

- In node A , calculate the (random) optimal decision in response to each value of d_1 :

$$A^*(d_1) = \arg \max_a \Psi_A(d_1, a).$$

This will provide us with the required distribution $p_D(A = a^* | d_1) = \mathbb{P}[A^*(d_1) = a]$, assuming that the space of a is discrete. It should be kept in mind again that the reduction of the previous influence diagram can be recast as

Algorithm 1: Overall attacker-defender approach.

```

1 Initialize parameters
2 For the Adversary
3 for each  $d_1$ , and for  $k = 1, \dots, n$  do
4   In node  $S_2$ , and for each  $(a, s_1, d_2)$ ,
5     Generate
6        $u_A^k(a, s_2) \sim U_A(a, s_2)$  for each  $(a, s_2)$ 
7        $p_A^k(S_2 = s_2 | d_2, a) \sim P_A(S_2 = s_2 | d_2, a)$  for  $(d_2, a)$ 
8     Calculate
9        $\psi_A^k(a, s_1, d_2) = \sum_{s_2} u_A^k(a, s_2) p_A^k(S_2 = s_2 | d_2, a)$ 
10    In node  $D_2$ , and for each  $(d_1, a, s_1)$ ,
11    Generate
12       $p_A^k(D_2 = d_2 | d_1, s_1) \sim P_A(D_2 = d_2 | d_1, s_1)$ 
13    In node  $S_1$ , and for each  $(d_1, a)$ ,
14    Generate
15       $p_A^k(S_1 = s_1 | d_1, a) \sim P_A(S_1 = s_1 | d_1, a)$ 
16    Calculate
17       $\psi_A^k(d_1, a) = \sum_{d_2} \psi_A^k(a, s_1, d_2) p_A^k(D_2 = d_2 | d_1, s_1)$ 
18    In node  $A$ ,
19    Calculate  $a_k^*(d_1) = \arg \max_a \psi_A^k(d_1, a)$ 
20  end
21 Approximate for each  $a$ ,
22    $\hat{p}_D(A = a | d_1) = |\{1 \leq k \leq n : a_k^* = a\}|/n$ 
23 For the Defender, calculate
24 In node  $S_2$ , for each  $(a, d_2)$ ,
25    $\psi_D(a, d_2) = \sum_{s_2} u_D(d_2, s_2) p_D(s_2 | d_2, a)$ 
26 In node  $D_2$ , for each  $(d_1, s_1)$ ,
27    $\psi_D(d_1, a, s_1) = \sum_{d_2} \psi_D(a, d_2) p_D(d_2 | d_1, s_1)$ 
28 In node  $S_1$ , for each  $(d_1, a)$ ,
29    $\psi_D(d_1, a) = \sum_{s_1} \psi_D(d_1, a, s_1) p_D(s_1 | d_1, a)$ 
30 In node  $A$ , for each  $d_1$ ,
31    $\psi_D(d_1) = \sum_a \psi_D(d_1, a) \hat{p}_D(a | d_1)$ 
32 In node  $D_1$ , calculate  $\psi_D = \arg \max_{d_1} \psi_D(d_1)$ .
```

$$A^*(d_1) = \arg \max_a \sum_{s_1} \sum_{d_2} \sum_{s_2} U_A(a, s_2) P_A(S_2 | d_2, a) \times \\ \times P_A(D_2 | d_1, s_1) P_A(S_1 | d_1, a).$$

The distribution $p_D(A|d_1)$ can be estimated by Monte Carlo simulation. To do this, we sample n times the proba-

bilities and utilities of the set

$$F = \{P_A(S_1|a, d_1), P_A(D_2|s_1, d_1), P(S_2|d_2, a), U(a, s_2)\}$$

to obtain the optimal decision $a^* \sim A^*(d_1)$ in the k -th iteration of the Monte Carlo simulation, $k = 1, \dots, n$. Then, we can approximate $p_D(A|d_1)$ through $|1 \leq k \leq n : a_k^* = a|/n$.

Note that, of the four components in F , the first three can be easily obtained. Normally, $P_A(S_1|a, d_1)$ would be related to $p_D(S_1|a, d_1)$ through some uncertainty, as we are referring to the results and the interaction between the attacker and defender, based on their decisions d_1 and a . This is also true for $P_A(S_2|d_2, a)$ with respect to $p_D(S_2|d_2, a)$. Regarding U_A , we will generally have information about the multiple interests of the adversary, which we add. However, the fourth element, $P_A(D_2|s_1, d_1)$, could require strategic thinking. In fact, the proposal presented here can be seen as a model of "level-2" thought, in which the defender optimizes their expected utility, with adverse probabilities derived from the optimization of the expected utilities (at random) of the adversary.

5.3 Overall Approach

The above ideas can be integrated into a step-by-step algorithm. First, we use simulation to estimate the distribution that predicts the options of the adversarial (suspicious) presence in the monitored sites. Second, we find the optimal initial allocation d_1^* (in favour or against the deployment of technology), maximizing the defender's expected utility with respect to the distribution of the derived prediction. We assume that the intervening r.v.'s are discrete, that is, that the impacts S_1 and S_2 are classified.

In the previous scheme, d_1 would be the optimal initial defence allocation. The corresponding probabilities of adversary presence, $p_D(A = a|d_1)$, represent the probability of each adversary scenario after deploying d_1^* , which would help raise awareness about the state of the security.

It is important to take into account that we make simulations for all possible initial allocations d_1 . Our problem is a binary allocation on using a technology or continuing with the status quo; however, in problems where the number of these allocations/resources is too large, we could use a regression meta-model, as explained in [21], simulating some defences, evaluating the corresponding attack probabilities and, consequently, approaching the attack probabilities in other defences. Then, we would use that estimated attack prediction distribution to find the optimal resource allocation.

6 EXPERIMENTAL EVALUATION

This section evaluates the decision-model proposed in the previous section. We use artificial data, given the absence of real, accurate information of terrorist Web-browsing data and counterterrorism strategies. Therefore, we give value to each of the evaluations that the defender must make about their own decisions and their beliefs about the adversary's decisions.

We also illustrate the proposed decision model with an example that will serve to show some of the computational subtleties, as well as the typical problem solving approach

used in a real case. In fact, it may serve as a template for real problems, which would basically add modelling and computational complexities. Essentially, first we structure the problem, and then model the defender's evaluations about themselves, and later, about the adversary. In the computational phase, we simulate the adversary's problem to obtain their attack probabilities and feed them into the defender's problem to obtain the optimal defence. Finally, we carry out a sensitivity analysis. For purposes of completeness and comparison, we also provide a standard game-theoretic approach under assumptions of common knowledge.

6.1 Structure of the Problem

We begin by identifying the available resources for both the defender and the adversary.

Defensive resources. We consider the defender's defensive resources to be the use of the automatic threat detection system with $d_1 = 1$. Otherwise, $d_1 = 0$.

Adversarial resources. For the adversary's resources we take the presence of the adversary in the set of monitored sites, with $a = 1$. Otherwise, $a = 0$.

Results of the game. Finally, we must consider the results of the decisions of both agents. We assume that the states of S_1 and S_2 will be 0 or 1, which means, respectively, the success or failure of the detection in terms of the alarm signal of the automatic system (recall that, if $d_1 = 0$, then $s_1 = 0$) and the final detection of the threat after a manual investigation.

6.2 The Defender's Evaluations

Now we consider the evaluations of the beliefs and preferences for the agency, that is, $p_D(S_1|a, d_1)$, $p_D(D_2|d_1, s_1)$, $p_D(S_2|d_2, a)$ and $u_D(d_2, s_2)$, defined in Sec. 5. In the assumed scenario, none of them will require strategic thinking. In the evaluations, we use the different parameters of the online surveillance problem defined in Sec. 3. Next, we examine the evaluation of the probability distributions involved as well as the utility model.

- *Evaluating $p_D(S_1|a, d_1)$.* S_1 represents the probability that the automatic threat detection system generates an alarm, whether there is suspicion or not. Obviously if the system is not used, it is impossible that it generates an alarm. We establish a range of values for this probability, although the defender will operate their problem with the base values.

	$a = 1$	$a = 0$
$d_1 = 1$	α^{base}	β^{base}
$d_1 = 0$	0	0

Table 3: $p_D(S_1 = 1|d_1, a)$.

- *Evaluating $p_D(D_2|d_1, s_1)$.* D_2 represents the probability of manually investigating a profile collected both when the automatic system is used and when it is not. We also establish a range of values that includes the base value for the defender's problem.
- *Evaluating $p_D(S_2|d_2, a)$.* S_2 represents the final success/failure of the surveillance. As we described in Sec. 3, manual investigation is considered 100% effective in confirming or ruling out a threat. In this case, we do not use a range for the values of this probability.

	$s_1 = 1$	$s_1 = 0$
$d_1 = 1$	ρ_1^{base}	ρ_0^{base}
$d_1 = 0$		ρ_0^{base}

Table 4: $p_D(D_2 = 1|d_1, s_1)$.

	$a = 1$	$a = 0$
$d_2 = 1$	1	0
$d_2 = 0$	0	0

Table 5: $p_D(S_2 = 1|a, d_2)$.

- *Evaluating $u_D(d_2, s_2)$.* Finally, the utility $u_D(d_2, s_2)$ as a measure of the quality of the model. We opted for an exponential utility function that allows us to order the costs u_D of the defender while assuming their (constant) risk aversion. Accordingly, we define $u_D(d_2, s_2) = -\exp(c_D v_D)$, with $c_D \sim U(0, 3)$.

	$s_2 = 1$	$s_2 = 0$
$d_2 = 1$	c	$c + d$
$d_2 = 0$	0	d

Table 6: $u_D(d_2, s_2)$.

6.3 The Defender's Evaluations about the Adversary

The security agency also needs to evaluate $p_D(A|d_1)$. This requires strategic thinking, as explained in Sec. 3. To do this, we must put ourselves in the adversary's shoes and make assessments about their probabilities and utilities, from the defender's perspective. Next, we go through how to estimate the probability distributions of the problem at hand and the adversary's utility function.

- *Evaluating $P_A(S_1|a, d_1)$.* We assume that $p_A(S_1 = 1|d_1, a)$ is similar to $p_D(S_1 = 1|d_1, a)$. To model our lack of knowledge about the probabilities used by the adversary in their decision problem, we add some uncertainty. In particular, we assume that, except in cases where $p_D(S_1 = 1|d_1, a)$ is 0 or 1, for those who suppose that the adversary's probabilities will match their beliefs P_A about $p_A(S_1 = 1|d_1, a)$ are uniform within the ranges $[p_A^{\min}, p_A^{\max}]$ of Table 7, evaluated by the defender.

	$a = 1$	$a = 0$
$d_1 = 1$	$[\alpha^{\min}, \alpha^{\max}]$	$[\beta^{\min}, \beta^{\max}]$
$d_1 = 0$	0	0

Table 7: $p_A(S_1 = 1|d_1, a)$.

Then we model p_A as a uniform distribution between p_A^{\min} and p_A^{\max} . Thus, $P_A(S_1|a, d_1)$, is defined by the expression

$$p_A = p_D^{\min} + \omega(p_D^{\max} - p_D^{\min}),$$

with $\omega \sim U(0, 1)$, so that the uncertainty about ω induces uncertainty about p_A to provide P_A .

- *Evaluating $P_A(D_2|d_1, s_1)$.* We adopt the same approach as before, now based on Table 8.
- *Evaluating $P_A(S_2|d_2, a)$.* We adopt the same approach as before, now based on Table 9.
- *Evaluating $U_A(a, s_2)$.* Finally, for utility $u_A(a, s_2)$ we also opted for an exponential utility function that allows us to order the adversary's costs v_A while we

	$s_1 = 1$	$s_1 = 0$
$d_1 = 1$	$[\rho_1^{\min}, \rho_1^{\max}]$	$[\rho_0^{\min}, \rho_0^{\max}]$
$d_1 = 0$	$[\rho^{\min}, \rho^{\max}]$	

Table 8: $p_D (D_2 = 1|d_1, s_1)$.

	$a = 1$	$a = 0$
$d_2 = 1$	1	0
$d_2 = 0$	0	0

Table 9: $p_D (S_2 = 1|a, d_2)$.

assume their (constant) risk seeking in relation to their benefits. Thus, we define $u_A(a, s_2) = \exp(c_A v_A)$, with $c_A \sim U(0, 0.025)$.

	$a = 1$	$a = 0$
$d_2 = 1$	$b - l$	b
$d_2 = 0$	0	0

Table 10: $v_A(a, s_2)$.

6.4 Results

We solved the problem with the open-source software R⁵ with an Intel® Core™ processor i3-2370 CPU at 2.4 GHz, 4Gb RAM on a Windows 10 64-bit operating system. In our example, the computation time is acceptable (between 15-20 seconds per problem on average) and therefore we will not consider the implementation and its performance as the object of the analysis. In any case, it should be noted that the resolution of the problem implies a Monte Carlo simulation for each value d_1 and that in each simulation we must propagate uncertainty at different levels, which becomes a strong computational challenge for larger problems.

For comparative and sensitivity analysis purposes, we prepared an experiment with 1000 random scenarios for five levels c_D of risk aversion of the defender, specifically $c_D = (0.01, 0.1, 0.5, 1.0, 3.0)$. In total, we obtained a set of 5000 solutions. In this way, we intend to determine how the parameters of the problem influence the defender's optimal decision d_1^* and at the same time compare their behaviour according to their level of risk aversion, between $c_D = 0.01$ (minimum) and $c_D = 3.0$ (maximum). In Table 11 we include the parameters that we used when the 1000 scenarios were generated for each value of c_D .

The implementation of our evaluation model allows us to obtain in each run the optimal solution d_1^* for the defender and the conditional probability $p_D(A|d_1)$ that the defender needed to solve their problem. As explained above, this probability is estimated using a Monte Carlo simulation with $n = 10000$ replications for each value $d_1 \in D$. Tables 12 and 13 show an extract of the results for one scenario and different levels of c_D of risk aversion of the defender, and the explicit form of the probability $p_D(A|d_1)$ for that scenario.

Thus, for example, the probability that the adversary is present in the set of monitored websites, taking into account the possibility that the defender is monitoring their navigation, is $\hat{p}_D(A = 1 | d_1 = 1) = 0.91$. This means that, in this example, solving the defender's problem ends with the optimal solution $d_1^* = 1$ for levels c_D of the defender's

5. R version 3.3.3 (2017-03-06).

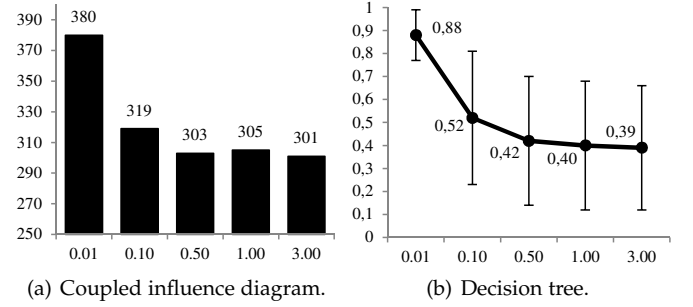


Figure 8: Results of ARA. No. of cases (over 1000) and ratio $\frac{\psi_D(d_1^{*ARA})}{\psi_D(d_1^{*})}$ (average \pm deviation) for d_1^* , both depending on the risk aversion level c_D of the defender (abscissa axis).

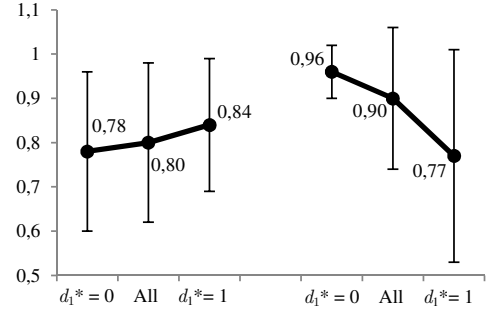


Figure 9: Values of $\hat{p}_D(A|d_1 = 1)$. On the left if $d_1 = 1$ and on the right if $d_1 = 0$.

risk aversion 0.01 and 0.10 and the contrary for higher levels.

In Fig. 8 and 9 we can observe graphically some of the most relevant results. The favourable use of the system, $d_1^* = 1$, is given in a moderate proportion of the 1000 cases (between 30% and 38%), and in a more conservative way at a higher level c_D of the defender's risk aversion. The same decreasing behaviour is observed for the ratio between the expected utilities of the optimal solution and its opposite. On the other hand, we have the average values of the estimates, $\hat{p}_D(A = 1 | d_1)$, for which we observe that $\hat{p}_D(A = 1 | d_1 = 1) \geq \hat{p}_D(A = 1 | d_1 = 0)$ when $d_1^* = 1$, and conversely in the opposite case. All these results confirm what we intuitively assumed a priori, and the correct behaviour of the calculations.

Finally, we adjust a parametric model to the set of solutions, specifically a logistic regression of the form $\text{logit}(d_1^*) = b_0 + b_1x_1 + b_1x_1 + \dots + b_ix_i$, with the aim of determining the relationship between d_1^* and the parameters of the problem. To select the best model we used the `bestglm`⁶ package of R, in an exercise set to avoid losing information and overestimating the logit model. Table 14 shows the results of the adjustments, where between the null model and the complete model, the best model obtained is 'ARA08.06' (logit model with 6 variables out of 8 available variables, highlighted in bold). Thus, the model indicates that, a priori, we could do without the parameters β^{base} and λ to explain the optimal decision $d_1^* = 1$ of the defender, while the parameter ρ^{base} (proportion of profiles investigated manually when the system is not used), with an *odds ratio* = 16.56, is shown to be highly influential.

6. <https://cran.r-project.org/web/packages/bestglm/index.html>

Table 11: Main parameters of the experimental evaluation.

Sensitivity and specificity		Proportion of manual investigations				Costs and coefficients		
α^{base}	β^{base}	ρ^{base}	ρ_1^{base}	ρ_0^{base}	ϕ	λ		
$U(0.60, 0.99)$	$U(0, 0.1)$	$U(0, 1)$			$U(0, 1)$	$U(0, 1)$		
α^{min}	β^{min}	ρ^{min}	ρ_1^{min}	ρ_0^{min}	c	b		
$U(0.60, \alpha^{\text{base}})$	$U(0, \beta^{\text{base}})$	$U(0, \rho^{\text{base}})$	ditto for ρ_1^{min} and ρ_0^{min}		$\phi d\phi \sim U(0, 1)$	100		
α^{max}	β^{max}	ρ^{max}	ρ_1^{max}	ρ_0^{max}	d	l		
$U(\alpha^{\text{base}}, 0.99)$	$U(\beta^{\text{base}}, 0.1)$	$U(\rho^{\text{base}}, 1)$	ditto for ρ_1^{max} and ρ_0^{max}		100	$l = (1 + \lambda)b\lambda \sim U(0, 1)$		

d_1^*	c_D	$\frac{\psi_D(d_1^*)}{\psi_D(d_1^{*-})}$	α^{base}	β^{base}	ρ^{base}	ρ_1^{base}	ρ_0^{base}	ϕ	λ	$\hat{p}_D(A = 1 d_1)$	
										$d_1 = 1$	$d_1 = 0$
1	0.01	0.73	0.92	0.01	0.48	0.89	0.73	0.08	0.05	0.91	0.99
1	0.10	0.53	0.92	0.01	0.48	0.89	0.73	0.08	0.05	0.91	0.99
0	0.50	0.21	0.92	0.01	0.48	0.89	0.73	0.08	0.05	0.91	0.99
0	1.00	0.08	0.92	0.01	0.48	0.89	0.73	0.08	0.05	0.91	0.99
0	3.00	0.08	0.92	0.01	0.48	0.89	0.73	0.08	0.05	0.91	0.99

Table 12: Summary of the results obtained for our ARA model.

	$a = 1$	$a = 0$
$d_1 = 1$	0.91	0.09
$d_1 = 0$	0.99	0.01

Table 13: Form of $\hat{p}_D(A|d_1)$.

The confusion matrix⁷ shown in Table 15, computed from the predictions of the 'ARA08.06' model, indicates that 3703 of the 5000 scenarios (74%) would be correctly predicted, 953 out of 1608 (59%) refer to the cases where $d_1^* = 1$.

At this point we want to highlight that other parametric and/or nonparametric adjustments can be used alternatively to logistic regression. It is also possible to analyse the sensitivity of the problem from other angles, be it for example through game theory in its classical form or through differential calculus, in order to find solutions and optimal parameter configurations. We will look at this in the following subsection, which analyses the problem from a standard game-theoretic approach.

6.4.1 Comparison with Game Theory

Now we will solve our example using the game theory approach and compare the results with the solution offered by ARA. The standard game theory approaches generally assume a common knowledge about the structure of the game (values at stake, resources available for the players and feasible assignments, etc.) as well as the utilities and probabilities of the players. In addition, the existence of objective probabilities for uncertainties is also usually assumed, in our case the results of automatic and/or manual investigations, $S_1|d_1, a$ and $S_2|a, d_2$.

We assume that the conditional probabilities ($S_1|d_1, a$) and $p(S_2|a, d_2)$ derive, respectively, from $p_D(S_1|d_1, a)$ and $p_D(S_2|a, d_2)$, that is, the defender's belief about the probability that the adversary's presence is not detected. These probabilities now represent objective non-detection probabilities, and both the defender and the adversary know them. In addition, the assumption of common knowledge ensures that the defender knows the probabilities used by the adversary when they solve their decision problem, and therefore does not need to represent the uncertainty surrounding them.

7. The cut-off value is 0.40 for a data set with 1608 reference cases over 5000.

To resolve the problem we adapted the reduction algorithm of influence diagrams proposed by [22] to evaluate an influence diagram that represents a decision problem of a single agent, to solve the sequential defence-attack games formulated as multi-agent influence diagrams. We solved, in parallel, the experiment proposed for ARA when the coefficients of risk aversion and risk seeking for the defender and the adversary are, respectively, $c_D = (0.01, 0.1, 0.5, 1.0, 3.0)$ and $c_A = 0.0125$ (remember that originally $c_A \sim U(0, 0.025)$). These values determine the utility functions of the defender and the adversary, given respectively by $u_D(d_2, s_2)$ y $u_A(a, s_2)$, that are also common knowledge. In general, we expect coincidental and opposite results that reveal the different assumptions of the methods. In Table 16, we show the same extract of results that Table 12 showed obtained by ARA, now adding the optimal solution provided by game theory. In this scenario and depending on the risk aversion level c_D of the defender, ARA behaves more prudently than game theory, which constantly obtains the same solution $d_1^{*GT} = 1$.

In Fig. 10, we can observe graphically some of the results obtained. The favourable use of the system, $d_1^* = 1$, is given in a proportion between 49% and 59% of 1,000 cases, decreasing to a greater level c_D of risk aversion of the defender. The same decreasing behaviour is observed for the ratio between the expected utilities of the optimal solution and its opposite. Compared to ARA, the frequency of favourable use of the system is significantly less conservative. This result satisfies us because it corresponds to ARA applications solving other problems. The opposite occurs with the ratio of the expected utilities between the optimal solution and the opposite, where ARA also has decreasing but higher ratios. We understand that these differences are found in the terms used to calculate the expected utilities of the defender in node A.

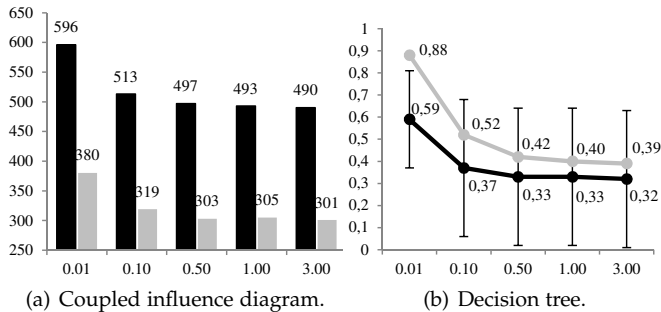
In Table 17, we show the same logistic adjustment used for the results obtained with ARA. Between the null model and the complete model, the best model obtained with game theory is 'GT08.05' (logit model with 5 variables out of 8 available variables). The model indicates that, a priori, we could dispense with parameters α^{base} , ϕ and λ to explain the defender's optimal decision, while parameter β^{base} (false positives), with an *odds ratio* of 29.88, demonstrates to be highly influential. For comparison, we also include in the table the best ARA model obtained, 'ARA08.06'.

Model	cte. logit	α^{base}	β^{base}	ρ^{base}	ρ_1^{base}	ρ_0^{base}	ϕ	λ	c_D	AIC
ARA00.00	-0.75	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	6283.70
ARA08.06	0.64	-1.49	0.00	2.81	-0.46	-2.37	-0.58	0.00	-0.08	5281.03
ARA08.08	0.59	-1.46	-1.08	2.81	-0.47	-2.37	-0.58	0.16	-0.08	5277.06

Table 14: Logit models for the ARA results.

		Pred.	
		$d_1^* = 0$	$d_1^* = 1$
Obs.	$d_1^* = 0$	2753	639
	$d_1^* = 1$	655	953

Table 15: Confusion matrix of the ‘ARA08.06’ logit model.

Figure 10: ARA results (grey) vs. game theory (black). No. of cases (over 1000) and ratio $\frac{\psi_D(d_1^{*GT})}{\psi_D(d_1^{*})}$ (average \pm deviation) for $d_1^* = 1$, both depending on the level c_D of risk aversion of the defender.

The confusion matrix⁸ shown in Table 18, computed from the predictions of the ‘GT08.05’ model, indicates that 3067 of the 5000 scenarios (61%, 74% with ARA) would be correctly predicted, 1312 out of 2589 (51%, 59% with ARA) refer to the cases in which $d_1^* = 1$.

In general, we can conclude that the results obtained with ARA are more satisfactory than those obtained with game theory. In contrast to standard game theory, however, ARA provides more costly solutions from a computational point of view, but more realistic in terms of the dynamics and perspective of the game proposed between adversaries to solve the problem of online surveillance.

7 DISCUSSION

We have defined the problem of online surveillance based on the intrusion-detection problem posed by [14] through backward induction. One of the main limitations of the mentioned work is the assumption of common knowledge of the parameters of the model by the agents. On this basis, we have investigated an ARA model as a novel way of proposing and solving the problem of online surveillance, where, among other aspects, we do not assume the hypothesis of common knowledge. Unlike the problem tackled by [14], we evaluate the problem of online surveillance faced by a security agency that monitors a set of specific websites by tracking and classifying profiles that are potentially suspected of carrying out terrorist attacks.

Our analysis constitutes a preliminary, theoretical step in that it aims to establish a point of departure and connection between the analytical framework provided by ARA, a young field within risk analysis, and the problem of online surveillance with counterterrorist purposes, understood as

8. The cut-off value is 0.55 for a data set with 2589 reference cases over 5000.

a game between opponents who want to maximize their benefits.

In order to give consistency to our proposal, we have illustrated a feasible architecture for online surveillance based on an engine for tracking user navigation traces on monitored websites and an automatic classifier of suspects, thought of as a classification method based on artificial intelligence. In this scenario, we have evaluated the adoption of automatic technology compared to the status quo that involves the manual investigation of profiles. We have applied the ARA methodology, which offers the possibility of treating the problem with game theory and risk analysis approaches in a new perspective of decision analysis against intelligent adversaries, who increase the risk of security and uncertain results.

Our experimental results corroborate the benefits of the proposed model and at the same time are indicative of its potentialities. Compared to the analysis of the problem from standard game theory, ARA indicates in general greater prudence in the deployment of the automatic system, even more the higher the level of risk aversion. In addition, the behavior of estimated conditional probabilities correctly responds to our intuitions ($\hat{p}_D(A = 1|d_1 = 1) \geq \hat{p}_D(A = 1|d_1 = 0)$ when $d_1^* = 1$, and conversely in the opposite case). From this point of view, the ARA model is more attractive than the standard game-theoretic model since it behaves satisfactorily without having to relax crucial hypotheses such as common knowledge and therefore subtracting realism from the problem.

We have used a parametric model with the aim of understanding the relationship between the optimal decision $d_1^* = 1$ and the parameters of the problem. This can be a way of not having to execute the resolution algorithm countless times to determine the adjustment of the system parameters. This and other decisions can be part of the implementation of the online monitoring architecture (machine learning, sensitivity analysis, etc.). In addition to the better parametric adjustment of ARA over standard game theory, we have observed that ρ^{base} (proportion of profiles investigated manually in case of not using the system) is determinant for ARA while β^{base} (false positives) is for game theory. β^{base} depends on how good our detection system is and ρ^{base} on the available resources. A priori the two implications could be coherent and a further debate could explore this, for example, depending on how the surveillance architecture is implemented.

We consider that the obtained results, although based on artificial and therefore limited data, place us positively at this starting point, as they have been satisfactorily contrasted with standard game theory. ARA goes beyond the dynamics of a player facing ‘‘nature’’, introducing in its place intelligent adversaries in a game of rational confrontation represented by the agents’ utility functions. In a nutshell, we can conclude that ARA is an excellent option for modeling and solving the problem posed against the classical model of game theory.

c_D	d_1^{*GT}	$\frac{\psi_D(d_1^{*GT})}{\psi_D(d_1^{*-})}$	d_1^{*ARA}	$\frac{\psi_D(d_1^{*ARA})}{\psi_D(d_1^{*-})}$	α^{base}	β^{base}	ρ^{base}	ρ_1^{base}	ρ_0^{base}	ϕ	λ
0.01	1	0.77	1	0.73	0.92	0.01	0.48	0.89	0.73	0.08	0.05
0.10	1	0.04	1	0.53	0.92	0.01	0.48	0.89	0.73	0.08	0.05
0.50	1	0.04	0	0.21	0.92	0.01	0.48	0.89	0.73	0.08	0.05
1.00	1	0.04	0	0.08	0.92	0.01	0.48	0.89	0.73	0.08	0.05
3.00	1	0.04	0	0.08	0.92	0.01	0.48	0.89	0.73	0.08	0.05

Table 16: Extract of standard game theory (GT) results vs. ARA results.

Model	Cte. logit	α^{base}	β^{base}	ρ^{base}	ρ_1^{base}	ρ_0^{base}	ϕ	λ	k_D	AIC
GT00.00	0.07	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	6927.13
GT08.05	0.34	0.00	3.40	1.34	-1.25	-0.74	0.00	0.00	-0.09	6549.20
GT08.08	-0.13	0.50	3.28	1.35	-1.24	-0.74	0.05	0.10	-0.09	6550.38
ARA08.06	0.64	-1.49	0.00	2.81	-0.46	-2.37	-0.58	0.00	-0.08	5281.03

Table 17: Logit models for the game theory results.

		Pred.	
		$d_1^* = 0$	$d_1^* = 1$
Obs.	$d_1^* = 0$	1595	816
	$d_1^* = 1$	1277	1312

Table 18: Confusion matrix of the 'GT08.05' logit model.

8 CONCLUSION AND FUTURE WORK

In recent years, Western countries have allocated tremendous amounts of resources to fight terrorism. As in any war, the battle occurs in various environments and the Internet, with the advent of the IoT, is one of the most powerful ones for propagating a threat and recruiting terrorists. However, at the very same time, this environment is the perfect storm for the development of ubiquitous online surveillance.

In this work, we have first examined the suitability of standard game theory and ARA, to tackle the online surveillance problem in which a security agency aims at countering terrorism online by deploying an automatic threat detection system on certain target websites. Then, we have proposed an ARA-based model to analyze the feasibility of using such an automatic system, and have determined under which conditions said deployment is better than the traditional model in which terrorist online activity is inspected manually by agents. Experimental results show that our ARA-based model is more attractive than the standard game-theoretic model as the former behaves satisfactorily without having to relax crucial hypotheses such as common knowledge and thus subtracting realism from the problem.

Future research lines include adopting other sequences and/or introducing new intermediate decisions to be taken into account (for example, changing the uncertainty node D2 to a decision node). For this, we can use other ARA templates and model new situations, both in sequential and simultaneous game dynamics.

ACKNOWLEDGMENTS AND DISCLAIMER

Partial support to this work has been received from the Spanish Government (project TIN2016-80250-R "SecMCloud"). J. Parra-Arnau is the recipient of a Juan de la Cierva postdoctoral fellowship, IJCI-2016-28239, from the Spanish Ministry of Economy and Competitiveness. J. Parra-Arnau is with the UNESCO Chair in Data Privacy, but the views in this paper are their own and are not necessarily shared by UNESCO.

REFERENCES

- [1] F. Reinares and C. G. Calvo, "Estado islámico en españa," Elcano Royal Inst., Tech. Rep., 2016.
- [2] A. Hintz and L. Denck, "The politics of surveillance policy: UK regulatory dynamics after snowden," *Inform. Sys. Res.*, vol. 5, no. 3, 2016.
- [3] V. Bier and M. Azaie, *Game Theoretic Risk Analysis of Security Threats*. Boston, US: Springer-Verlag, 2008, vol. 128.
- [4] I. Rios, J. Rios, and D. Banks, "Adversarial risk analysis," *J. Amer. Stat. Assoc.*, vol. 104, no. 486, pp. 841–854, 2009.
- [5] J. Parra-Arnau and C. Castelluccia, "On the cost-effectiveness of mass surveillance," *IEEE Access*, vol. 6, no. 1, pp. 46 538–46 557, Dec. 2018.
- [6] S. Englehardt and A. Narayanan, "Online tracking: A 1-million-site measurement and analysis," in *Proc. ACM Conf. Comput., Commun. Secur. (CCS)*. ACM, 2016, pp. 1388–1401.
- [7] J. Parra-Arnau, J. P. Achara, and C. Castelluccia, "MyAdChoices: Bringing transparency and control to online advertising," *ACM Trans. Web*, vol. 11, no. 1, Mar. 2017. [Online]. Available: <https://hal.inria.fr/hal-01270186/document>
- [8] S. Yuan, A. Z. Abidin, M. Sloan, and J. Wang, "Internet advertising: An interplay among advertisers, online publishers, ad exchanges and web users," *arXiv: 1206.1754*, 2012, arXiv preprint.
- [9] V. Toubiana, "SquiggleSR," 2007. [Online]. Available: www.squigglesr.com
- [10] B. Liu, A. Sheth, U. Weinsberg, J. Chandrashekar, and R. Govindan, "Adreveal: Improving transparency into online targeted advertising," in *Proc. Hot Topics in Netw.* ACM, 2013, pp. 121–127.
- [11] J. Yan, N. Liu, G. Wang, W. Zhang, Y. Jiang, and Z. Chen, "How much can behavioral targeting help online advertising?" in *Proc. Int. WWW Conf.* ACM, 2009, pp. 261–270.
- [12] M. Aly, A. Hatch, V. Josifovski, and V. K. Narayanan, "Web-scale user modeling for targeting," in *Proc. Int. WWW Conf.* ACM, 2012, pp. 3–12.
- [13] M. M. Tsang, S. C. Ho, and T. P. Liang, "Consumer attitudes toward mobile advertising: An empirical study," *Int. J. Electron. Commer.*, vol. 8, no. 3, pp. 65–78, 2004.
- [14] H. Cavusoglu, B. Mishra, and S. Raghunathan, "The value of intrusion detection systems in information technology security architecture," *Inform. Sys. Res.*, vol. 16, no. 1, pp. 28–46, 2005.
- [15] P. Mell and R. Bace, "Nist special publication on intrusion detection systems," National Institute of Standards and Technology (NIST), Tech. Rep., 2001.
- [16] J. Merrick and L. McLay, "Is screening cargo containers for smuggled nuclear threats worthwhile?" *Decision Anal.*, vol. 7, no. 2, pp. 155–171, 2010.
- [17] J. Rios and D. R. Insua, "Adversarial risk analysis: Applications to basic counterterrorism models," in *Proc. Int. Conf. Alg. Decision Theory*, ser. Lecture Notes Comput. Sci. (LNCS), 2009, pp. 306–315.
- [18] J. Zhuang and V. Bier, "Balancing terrorism and natural disasters—defensive strategy with endogenous attacker effort," *Oper. Res.*, vol. 55, no. 5, pp. 976–991, Sep. 2007.
- [19] G. Brown, M. Carlyle, J. Salmeron, and K. Wood, "Defending critical infrastructure," *Interf.*, vol. 36, no. 6, pp. 530–544, Sep. 2006.
- [20] D. Winterfeldt and T. O'Sullivan, "Should we protect commercial airplanes against surface-to-air missile attacks by terrorists?" *Decision Anal.*, vol. 3, no. 2, pp. 63–75, 2006.
- [21] R. Barton and M. Meckesheimer, *Handbooks in Operations Research and Management Science*. Elsevier, 2006, ch. Metamodel-Based Simulation Optimization, pp. 535–574.

- [22] R. D. Shachter, "Evaluating influence diagrams," *Oper. Res.*, vol. 34, no. 6, 1996.