# Reward-Based Deception with Cognitive Bias

Bo Wu, Murat Cubuktepe, Suda Bharadwaj, and Ufuk Topcu

*Abstract*—Deception plays a key role in adversarial or strategic interactions for the purpose of self-defence and survival. This paper introduces a general framework and solution to address deception. Most existing approaches for deception consider obfuscating crucial information to rational adversaries with abundant memory and computation resources. In this paper, we consider deceiving adversaries with bounded rationality and in terms of expected rewards. This problem is commonly encountered in many applications especially involving human adversaries. Leveraging the cognitive bias of humans in reward evaluation under stochastic outcomes, we introduce a framework to optimally assign resources of a limited quantity to optimally defend against human adversaries. Modeling such cognitive biases follows the so-called prospect theory from behavioral psychology literature. Then we formulate the resource allocation problem as a signomial program to minimize the defender's cost in an environment modeled as a Markov decision process. We use police patrol hour assignment as an illustrative example and provide detailed simulation results based on real-world data.

## I. Introduction

Deception refers to a deliberate attempt to mislead or confuse adversaries so that they may take strategies that are in the defender's favor [1]. Deception can limit the effectiveness of an adversary's attack, waste adversary's resources and prevent the leakage of critical information [2]. It is a widely observed behavior in nature for self-defence and survival. Deception also plays a key role in many aspects of human society, such as economics [3], warfare [4], game [5], cyber security [2] and so on.

In this paper, we focus on the scenario in which the adversary acts in an environment where this interaction is modeled as a Markov decision process (MDP) [6]. The adversary's aim is to collect rewards at each state of the MDP and the defender tries to minimize the accumulated reward through *deception*. Many existing approaches for deception rely on a rational adversary with sufficient memory and computation power to find its optimal policy [7], [1]. However, deceiving an adversary with only bounded rationality [8], i.e., one whose decisions may follow certain rules that deviate from the optimal action [9], has not been adequately studied so far. Deceiving an adversary with bounded rationality finds, for example, its application in intrusion detection and protection [10] or public safety [11]. Different from obfuscating sensitive system information to the adversary [12], [13], by deception, we mean that the defender optimally assigns a

limited resource to each state, such that the expected cost from defender's perspective (or equivalently, the reward for the adversary) incurred by an adversary can be minimized, even though the adversary is expecting more based on his cognitively biased view of rewards.

To deceive a human more effectively, it is essential to understand the human's cognitive characteristics and what affects his decisions (particularly with stochastic outcomes). Works in behavior psychology, e.g. [14], suggested that humans' decision-making follows intuition and bounded rationality. Empirical evidence has shown that humans tend to evaluate gains and losses differently in decision-making [15]. Humans tend to over-estimate the likelihood of low-probability events and underestimate the likelihood of high-probability events in a nonlinear fashion [16], [15]. Risk-sensitive measures, such as those in the so-called prospect theory [16], capture such biases and are widely used in psychology and economics to characterize human preferences. Furthermore, humans tend to make decisions that are often sub-optimal [17]. It is generally believed that such sub-optimality is the result of intuitive decisions or preferences that happen automatically and quickly without much reflection [17], [14]. Human decisions are subject to stochasticity due to the limited computational capacity and inherent noise [18]. Consequently, human decisions are often cognitively biased (have a different reward mechanism), probabilistic (have a stochastic action selection policy) and memoryless (only depends on the current state). These are the very characteristics of human decision-making we expect to account for in reward-based deception.

This paper investigates how one can deceive a human adversary by optimally allocating limited resources to minimize his rewards. We model the environment as an MDP to capture the choices available to a human decision-maker and their probabilistic outcomes. We consider opportunistic human adversaries, i.e., they usually do not have significant planning and only act based on immediately available rewards [11]. We describe the human adversary's policy to select different actions following the prospect theory and bounded rationality [8]. We model both the adversary's perceived reward and defender's cost (equivalently, the adversary's reward from the defender's point of view) as functions of the resources available at each state of the MDP. Additionally, we define a subset of the states in the MDP as sensitive states that the human adversary should be kept from visiting.

We then formulate the optimal resource allocation problem as a signomial program (SP) to minimize the defender's cost. SPs are a special form of nonlinear programming problems, and they are generally nonconvex. Solving non-

Bo Wu, Murat Cubuktepe, Suda Bharadwaj and Ufuk Topcu are with the Department of Aerospace Engineering and Engineering Mechanics, and the Oden Institute for Computational Engineering and Sciences, University of Texas, Austin, 201 E 24th St, Austin, TX 78712. email: {bwu3, mcubuktepe, suda.b, utopcu}@utexas.edu

convex NLPs is NP-hard [19] in general, and a globally optimal solution of an SP cannot be computed efficiently. SPs generalize geometric programs (GP), which can be transformed into convex optimization problems and then can be solved efficiently [20]. In this paper, we approximate the proposed SP to a GP. In numerical experiments, we show that this approach obtain locally optimal solutions of the SP efficiently by solving a number of GPs. We demonstrate the approach with a problem on the assignment of police patrol hour against opportunistic criminals [11].

The problem we study is closely related to the Stackelberg security game (SSG) which consists of an attacker and a defender that interact with each other. In SSG, the defender acts first with limited resources and then the attackers play in response [21]. SSG is a popular formalism to study security problems against human adversaries. Early efforts focused on one-shot games where an adversary can only take one move [22] without considering human's bounded rationality. Then repeated SSG was considered in wildlife security [10] and fisheries [23] where the defender and the adversary can have repeated interaction. However, neither of these papers considered how a human perceives probabilities, where the existence of nonlinear probability weighting curves is a well-known result in prospect theory [16]. Such phenomenon was taken into account in [24] and [25]. But [24] only studied one-shot games and [25] did not consider the adversaries may move from place to place.

The rest of this paper is organized as the following. We first provide the necessary preliminaries for stochastic environment modeling, human cognitive biases and decision-making in Section II. Then we formulate the human deception problem in terms of resource allocation in Section III and show that it can be transformed into a signomial program in Section IV. We propose the computational approach to solve the signomial program in Section V. Section VI shows simulations results and discusses their implications. We conclude our paper and discusses possible future directions in Section VII.

## II. PRELIMINARIES

### A. Monomials, Posynomials, and Signomials.

Let $V = \{x_1, \ldots, x_n\}$ be a finite set of strictly positive real-valued *variables*. A *monomial* over $V$ is an expression of the form

$$f = c \cdot x_1^{a_1} \cdots x_n^{a_n} \ ,$$

where $c \in \mathbb{R}^+$ is a positive coefficient, and $a_i \in \mathbb{R}$ are exponents for $1 \leq i \leq n$. A *posynomial* over $V$ is a sum of one or more monomials:

$$g = \sum_{k=1}^{K} c_k \cdot x_1^{a_{1k}} \cdots x_n^{a_{nk}} \ . \tag{1}$$

If $c_k$ is allowed to be a negative real number for any $1 \leq k \leq K$, then the expression (1) is a *signomial*.

This definition of monomials differs from the standard algebraic definition where exponents are positive integers with no restriction on the coefficient sign. A sum of monomials is then called a *polynomial*.

### B. Nonlinear programs.

A general nonlinear program (NLP) over a set of real-valued variables $V$ is

$$\text{minimize} \quad f \tag{2}$$
$$\text{subject to}$$
$$g_i \leq 1, \quad i = 1, \ldots, m, \tag{3}$$
$$h_j = 1, \quad j = 1, \ldots, m, \tag{4}$$

where $f$, $g_i$, and $h_j$ are arbitrary functions over $V$, and $m$ and $p$ are the number of inequality and equality constraints of the program respectively.

### C. Signomial programs and geometric programs.

A special class of NLPs known as *signomial programs* (SP) is of the form (2)–(4) where $f$, $g_i$ and $h_j$ are signomials over $V$, see Def. II-A. A geometric program (GP) is an SP of the form (2)–(4) where $f, g_i$ are posynomial functions and $h_j$ are monomial functions over $V$. GPs can be transformed into convex programs [20, §2.5] and then can be solved efficiently using interior-point methods [26]. SPs are non-convex programs in general, and therefore there is no efficient algorithm to compute global optimal solutions for SPs . However, we can efficiently obtain local optimal solutions for SPs in our setting, as shown in the following sections.

In this paper, the adversary with bounded rationality moves in an environment modeled as a Markov decision process (MDP) [6].

### D. Markov Decision Processes.

A (MDP) is a tuple $\mathcal{M} = (S, \nu, A, T, U)$ where
- $S$ is a finite set of states;
- $\nu : S \rightarrow [0, 1]$ is the initial state distribution;
- $A$ is a finite set of actions;
- $T(s, a, s') := P(s'|s, a)$. That is, the probability of transiting from $s$ to $s'$ with action $a$; and
- $U(s) \in \mathbb{R}^+$ is the utility function that assigns resources with a quantity $U(s)$ to state $s$.

At each state $s$, an adversary has a set of actions available to choose. Then the nondeterminism of the action selection has to be resolved by a policy $\pi$ executed by the adversary. A (memoryless) policy $\pi : S \times A \rightarrow [0, 1]$ of an MDP $\mathcal{M}$ is a function that maps every state action pair $(s, a)$ where $s \in S$ and $a \in A$ with probability $\pi(s, a)$.

By definition, the policy $\pi$ specifies the probability for the next action $a$ to be taken at the current state $s$. A bounded rational adversary is often limited in memory and computation power, therefore we only consider the memoryless policies.

In an MDP, a finite state-action path is $\omega = s_0 a_0 s_1 a_1 \ldots$, where $s_i \in S, a_i \in A$ and $T(s_i, a_i, s_{i+1}) > 0$. Given a policy $\pi$, it is possible to calculate the probability of such path $P_\pi(\omega)$ as

$$P_\pi(\omega) = \nu(s_0) \prod_i \pi(s_i, a_i) T(s_i, a_i, s_{i+1}). \tag{5}$$

## III. REWARD-BASED DECEPTION

We assume that an adversary with bounded rationality moves around in an environment modeled as an MDP $\mathcal{M} = (S, \nu, A, T, U)$. When the adversary is at a state $s \in S$, from the defender's point of view, the immediate reward for the human adversary (or equivalently, the cost for the defender) is

$$R(s) = g(U(s)) \in \mathbb{R}^+,$$

which is a function of allocated resource $U(s)$. However, due to the bounded rationality and cognitive biases, the perceived immediate reward $R_h(s)$ at state $s$ by the adversary is a different function of $U(s)$, and is given by

$$R_h(s) = f(U(s)) \in \mathbb{R}^+,$$

where $f$ is another function over $U$. For a given policy $\pi$, expected rewards $Q_t(s)$ at each state $s$ and time t with a finite time horizon $H$ can be evaluated as

$$Q_t(s) = R(s) + \sum_a \sum_{s'} \pi(s, a) T(s, a, s') Q_{t+1}(s'), \quad (6)$$

where $t = 0..., H - 1$, $Q_H(s) = R(s)$. Therefore, $Q_t$ represents the expect accumulated cost of the defender, or equivalently, expected rewards for the human adversary obtained from the policy $\pi$.

The defender's objective is to optimally assign the resources to each state to minimize his cost (equivalently, the adversary's reward) $Q$, where

$$Q = \sum \nu(s) Q_0(s), \quad (7)$$

by designing the utility function $U$, where the resources are of limited quantity, i.e., $\sum_s U(s) = D$. Also imagine that there are set of sensitive states $S_s \subset S$ that the adversaries should be kept away from. Denote the set of paths that reach $S_s$ in $H$ steps as $\Omega$ such that for each $\omega \in \Omega$ where $\omega = s_0 a_0, ..., s_N$, we require $N \leq H, s_i \notin S_s, i < N$ and $s_N \in S_s$. In particular, given a policy $\pi$, $P(\lozenge^{\leq H} S_s)$ can be calculated as

$$P(\lozenge^{\leq H} S_s) = \sum_{\omega \in \Omega} P(\omega). \quad (8)$$

*Problem 1:* Given an MDP $\mathcal{M} = (S, \nu, A, T, U)$, time horizon $H$, human reward function defined in (12) and the policy $\pi$, design reward function $U$ with a limited total budget $\sum_s U(s) = D$, such that $Q$ as defined in (7) can be minimized $\hat{s}$ and $P(\lozenge^{\leq H} S_s) \leq \lambda$, which requires that the probability to reach $S_s$ in $H$ steps should be no larger than $\lambda \in [0, 1]$, i.e.,

$$P(\lozenge^{\leq H} S_s) \leq \lambda. \quad (9)$$

*Remark 1:* Problem 1 studies how to optimally assign the reward to trick the adversary into thinking that his policy could obtain more rewards but in fact, the actual expected reward is minimized with a low probability of visiting sensitive states $S_s$.

### A. Human Adversaries with Cognitive Biases

To solve Problem 1, it is essential to find the adversary's policy $\pi$. In this paper, we take human as the adversary with bounded rationality who is opportunistic, meaning that he does not have a specific attack goal nor plans strategically, but is flexible about his movement plan and seek opportunities for attacks [27]. Those attacks may incur rewards to the human adversary and consequently certain costs for the defender. The process of human decision-making typically follows several steps [28]. First, a human recognizes his current situation or state. Second, he will evaluate each available action based on the potential immediate reward it can bring. Third, he will select an action following some rules. Then he will receive a reward and observe a new state. In this section, we will introduce the modeling framework for the second and third step.

For a human with bounded rationality, the value of a reward from an action is a function of the possible outcomes and their associated probabilities. The prospect theory developed by Kahneman and Tversky [16] is a frequently used modeling framework to characterize the reward perceived by a human. Prospect theory claims that humans tend to overestimate the low probabilities and underestimate the high probabilities in a nonlinear fashion. For example, between winning 100 dollar with $\frac{1}{100}$ probability and nothing else, or 1 dollar with probability 1, humans tend to prefer the former, even though both have the same expectation.

Given $X$ as the discrete random variable that has a finite set of outcomes $O$, a general form of prospect theory utility $V(X)$ (i.e. the reward anticipated by a human) is the following.

$$V(X) = \sum_{x \in O} v(x) w(p(x)), \quad (10)$$

where $v(x) \in \mathbb{R}$ denotes the reward perceived by a human from the outcome $x$. The probability $p(x)$ to get the outcome $x$ is weighted by a nonlinear function $w$ that captures the human tendency to over-estimate low probabilities and under-estimate high probabilities.

The expected immediate reward $r_a(s)$ to perform an action $a$ at state $s$ is

$$r_a(s) = \sum_{s'} R(s') T(s, a, s'). \quad (11)$$

However, according to prospect theory, from a human's perspective, the perceived expected immediate reward $r_a^h(s)$ is different. Let $X_{s,a}$ be the random variable for the outcome $O_{s,a}$ of executing action $a$ at state $s$. We have $O_{s,a} = \{x_s' | T(s, a, s') > 0\}$ where $x_{s'}$ denotes the event that the state transits from $s$ to $s'$ with an action $a$. The distribution of $X_{s,a}$ is defined as follows.

$$p(x_{s'}) = T(s, a, s'), \forall x_{s'} \in O_{s,a}.$$

The human perceived reward $v(x_{s'})$ for the outcome $x_{s'}$ depends on $U(s')$ received from reaching the state $s'$, which is denoted by
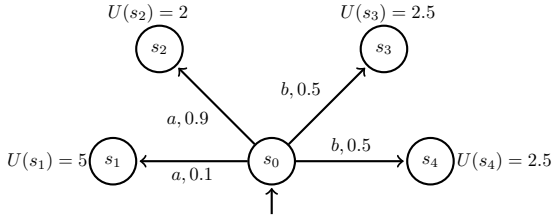
$$v(x_{s'}) = R_h(s') = f(U(s')).$$

Fig. 1. A simple example for sub-optimality with human cognitive biases

As a result, $r_a^h(s)$ is denoted by

$$r_a^h(s) = \sum_{x_{s'} \in O_{s,a}} v(x_{s'})w(p(x_{s'}))$$
$$= \sum_{s'} f(U(s'))w(T(s,a,s')). \quad (12)$$

An empirical form of $w$ is the following [16].

$$w(p) = \frac{p^\gamma}{(p^\gamma + (1-p)^\gamma)^{\frac{1}{\gamma}}}, \gamma > 0. \quad (13)$$

Given an MDP as depicted in Figure 1, where $S = \{s_0, \ldots, s_4\}$, $A = \{a, b\}$. We assume that $R(s) = U(s)$, $R_h(s) = U(s)^{0.88}$, $\gamma = 0.6$ in (13). It can be found from (11) and (12) that $r_a(s_0) = 2.3, r_a^h(s_0) = 2.0678, r_b(s_0) = 2.5$ and $r_b^h(s_0) = 1.8617$. Since $r_a^h(s_0) > r_b^h(s_0)$, suppose a human is at $s_0$, from human's perspective, he will prefer the action $a$. However, $r_a(s_0) < r_b(s_0)$ which indicates that action $a$ actually has more expected immediate rewards.

*Remark 2:* In this example, the rewards are already given, and it can be seen that the human could make a sub-optimal decision. It illustrates how cognitive bias can deviate the human behavior from optimal.

After evaluating the outcome of each candidate action $a$ by $r_a^h(s)$, a human then needs to make an action selection. Humans are known to only have quite limited cognitive capabilities. Human's policy $\pi$ to choose an action can be described as a random process that biases toward the actions of high $r_a^h(s)$, such that

$$\pi(s,a) = \frac{r_a^h(s)}{\sum_{a'} r_{a'}^h(s)}, \quad (14)$$

where $\pi(s,a)$ denotes the probability of executing the action $a$ at state $s$. Such a bounded rational behavior has been observed in humans, such as urban criminal activities [29]. Intuitively, it implies that human selects the action $a$ opportunistically at each state $s$ with the probability proportional to the perceived immediate reward $r_a^h(s)$.

Now we are ready to redefine Problem 1 as follows.

*Problem 2:* Solve Problem 1 for $\pi$ defined as (14).

### IV. SIGNOMIAL PROGRAMMING FORMULATION

Given an MDP $\mathcal{M}$, time horizon $H$, human reward function and policy as defined in (12) and (14), the solution of the Problem 1 can be computed by solving the following signomial program. The $g$ and $f$ are assumed to be monomial functions of $U$ for our solution method.

$$\text{minimize} \quad Q = \sum \nu(s)Q_0(s) \quad (15)$$

subject to

$$\forall s \in S, t \in \{0, \ldots, H-1\},$$
$$Q_t(s) \geq R(s) + \sum_{a \in A} \sum_{s' \in S} \pi(s,a)T(s,a,s')Q_{t+1}(s') \quad (16)$$

$$\forall s \in S, t \in \{0, \ldots, H-1\},$$
$$P_t(s) \geq \sum_{a \in A} \sum_{s' \in S} \pi(s,a)T(s,a,s')P_{t+1}(s') \quad (17)$$

$$\forall t \in \{0, \ldots, H\}, \forall s \in S_s, \quad P_t(s) = 1 \quad (18)$$

$$\forall t \in \{0, \ldots, H\}, \quad \sum_{s \in S} \nu(s)P_t(s) \leq \lambda \quad (19)$$

$$\forall s \in S, a \in A,$$
$$\pi(s,a) \sum_{a' \in A} \sum_{s' \in S} f(U(s'))w(T(s,a',s'))$$
$$= \sum_{s' \in S} f(U(s'))w(T(s,a,s')) \quad (20)$$

$$\forall s \in S, \quad R(s) = g(U(s)) \quad (21)$$

$$\sum_{s \in S} U(s) = D, \quad (22)$$

where variables $R(s)$ are for rewards in each state $s$, $U(s)$ are for utilities in each state $s$, $\pi(s,a)$ are for the probability of taking action $a$ in state $s$ are for each state and action, $Q_t(s)$ are for the expected reward of the state $s$ and time step $t$, and $P_t(s)$ are for the probability of reaching the set of target states $S_s$ in each state $s$ and time step $t$.

The objective in (15) minimizes the accumulated expected reward from the initial state distribution $\nu(s)$ over a time horizon $H$. In (16), we compute $Q_t(s)$ by adding the immediate reward in state $s$ and the expected reward of the successor states according to the policy variables $\pi(s,a)$ for each action $a$. The probability of reaching each successor state $s'$ depends on the policy variables $\pi(s,a)$ in each state $s$ and action $a$. Similar to the constraint in (16), the variables $P_t(s)$ are assigned to the probability of reaching the set of target states $S_s$ from state $s$ and time step $t$ in (17).

The probability of reaching any state $s \in S_s$ in each horizon from the states in $S_s$ is set to 1 as in (18). The constraint in (19) assures that the probability of reaching any state $s \in S_s$ from the initial state distribution $\nu(s)$ is less than $\lambda$. The constraint in (20) computes the policy using the model in (14). We give the relationship between rewards and utilities in (21). Finally, (22) gives the total budget for utilities.

The constraint in (16) and (17) are convex constraints, because the functions in the right hand sides are posynomial functions, and the functions in the left hand sides are monomial functions. The constraints in (18) and (19) are affine constraints, therefore they are convex. The constraints in (20) and (22) are equality constraints with posynomials, therefore they belong to the class of signomial constraints, and they are not convex. In the literature, there are various methods

to deal with the nonconvex constraints to obtain a locally optimal solution including sequential convex programming, convex-concave programming, branch and bound or cutting plane methods [30], [20], [31].

## V. COMPUTATIONAL APPROACH FOR THE SIGNOMIAL PROGRAM

In this section, we discuss how to compute a locally optimal solution efficiently for Problem 1 by solving the signomial program in (15)–(22). We propose a *sequential convex programming* method to compute a local optimum of the signomial program in (15)–(22), following [20, §9.1], solving a sequence of GPs. We obtain each GP by replacing signomial constraints in equality constraints of the SGP signomial program in (15)–(22) with *monomial approximations* of the functions.

### A. Monomial approximation

Given a posynomial $f$, a set of variables $\{x_1, \ldots, x_n\}$, and an initial point $\hat{x}$, a *monomial approximation* [20] $\hat{f}$ for $f$ around $\hat{x}$ is

$$\forall i.1 \leq i \leq n \quad \hat{f} = f[\hat{x}] \prod_{i=1}^{n} \left( \frac{x_i}{\hat{x}(x_i)} \right)^{a_i},$$

$$\text{where } a_i = \frac{\hat{x}(x_i)}{f[\hat{x}]} \frac{\partial f}{\partial x_i}[\hat{x}].$$

Intuitively, a monomial approximation of a posynomial $f$ around an initial point $\hat{x}$ corresponds to an affine approximation of the posynomial $f$. Such an approximation is provided by the first order Taylor approximation of $f$, see [20, §9.1] for more details.

For a given instantiation of the utility and policy variables $U(s)$ and $\pi(s, a)$, we approximate the SP in (15)–(22) to obtain a GP as follows. We first normalize the utility values to ensure that they sum up to $D$. Then, using those utility values, we compute the policy according to constraint in (20). After the policy comptutation, we compute a monomial approximation of each posynomial term in the constraints (20) and (22) around the previous instantiation of the utility and policy variables. After the approximation, we solve the approximate GP. We repeat this procedure until the procedure converges.

One key problem with this approach is, we require an initial feasible point to the signomial problem in (15)–(22), which may be hard to find because of the reachability constraint in (19). Therefore, we introduce a new variable $\tau$ and we replace the reachability constraint in (19) by the following constraints:

$$\forall t \in \{0, \ldots, H\}, \quad \sum_{s \in S} \nu(s) P_t(s) \leq \lambda \cdot \tau \quad (23)$$

$$\tau \geq 1. \quad (24)$$

By replacing the reachability constraint, we ensure that any initial utility function and policy is feasible to the signomial program in (15)–(24). To enforce the feasability of

the reachability constraint in (19), we change the objective in (15) as follows:

$$\text{minimize} \quad Q + \delta \cdot \tau \quad (25)$$

where $\delta$ is a positive *penalty parameter* that determines the violation rate for the *soft* constraint in (23). In our formulation, we increase $\delta$ after each iteration to satisfy the reachability constraint.

We stop the iterations when the change in the value of $Q$ is less than a small positive constant $\epsilon$. Intuitively, $\epsilon$ defines the required improvement on the objective value for each iteration; once there is not enough improvement, the process terminates.

## VI. NUMERICAL EXPERIMENT

Let us consider an urban security problem, where a criminal plans his next move randomly based on his local information on the nearby locations that are protected by police patrols. Such a criminal is opportunistic, i.e, he is not highly strategic by conducting careful surveillance and rational planning before making moves. It is known that this kind of opportunistic adversaries contribute to the majority of the urban crimes [32]. For prevention and protection, each location should be assigned a certain police patrol hours. Due to the limited amount of police resources, the total number of patrol hours is limited as well.
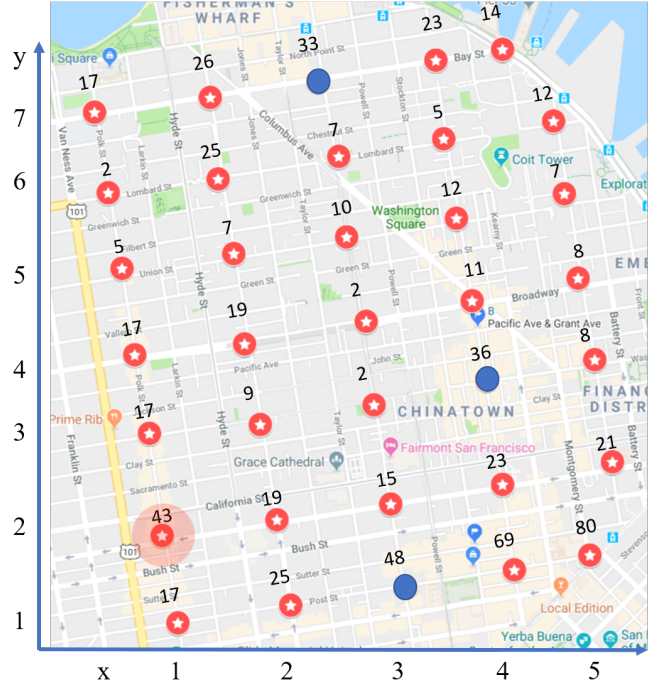


Fig. 2. The 35 intersections in North East of San Francisco. The map is obtained from Google map. The shaded area is a circle with a 500 feet radius. The numbering of the states starts from the bottom left corner and goes from left to right in every row. The number beside each location indicates the number of crimes in that area.

Figure 2 shows 35 intersections in San Francisco, CA with 7 rows and 5 columns. We use an MDP $\mathcal{M} = (S, \nu, A, T, U)$

to describe the network of the set of intersections $S$. The number $C(s)$ of crimes that occurred in the first four weeks of October, 2018 within 500 feet of each interaction $s$ is shown in Figure 2. The crime data are obtained from https://www.crimemapping.com/map/ca/sanfrancisco. The criminal can choose to move left, right, up or down to the immediate neighboring intersections. Consequently, there are four actions available. The execution of each action will lead the human to its intended neighborhood of the intersection with a high probability ($\geq 0.95$) and small probability to other neighboring intersections to account for unexpected change of movement plan.

Initially, the criminal has equal probability to appear at any state, i.e., $\nu(s) = \frac{1}{35}$ for any $s \in S$. The utility $U(s)$ denotes the number of police patrol hours that should be allocated to the vicinity of each intersection. The total number of police patrol hours is $D = \sum U(s) = 400$. If a location $s$ is assigned with $U(s)$ patrol hours, its reward to the criminal (equivalently, the cost to the defender) is

$$R(s) = \frac{C(s)}{U(s)}.$$

Intuitively, it means that the reward to the human adversary, from the defender's point of view, is proportional to the crime rate indicated by $C(s)$ and inversely proportional to the police patrol hours. The reward from the human adversary's view is evaluated as

$$f(U(s)) = R(s)^{0.88},$$

which is a function commonly seen in the literature to describe how human biases the reward [15].

Initially, the criminal is at $s$ with probability $\nu(s)$, where he tries to plan his move over the next $H$ steps. The objective is to assign the police patrol hours to each state, such that the expected accumulated reward in $H$ steps received by the criminal is minimized. The sensitive states $S_s = \{3, 14, 33\}$ should be visited with a probability no larger than $\lambda = 0.3$, i.e.

$$P(\lozenge^{\leq H} S_s) \leq 0.3.$$

The sensitive states are also shown as blue circles in Figure 2.

We formulate the problem as a signomial program. From an initial uniform utility distribution, we instantiate the policies and reward functions. Then, from the initial values, we linearize the signomial program in (15)–(25) to a geometric program. We parse the geometric programs using the tool GPkit [33], and solve them using the solver MOSEK [34]. We set $\epsilon = 10^{-4}$ for convergence tolerance. All experiments were run on a 2.3 GHz machine with 16 GB RAM. The procedure converged after 32 iterations for a problem with horizon length $T = 20$ in 230.06 seconds. The expected reward $Q$ from the initial state distribution is 117.15, and the reachability probability of the sensitive states from the initial state distribution is 0.192, which satisfies the reachability specification.

The result is shown in Figure 3. Different colors at each intersection show the number of patrol hours, i.e, the resource
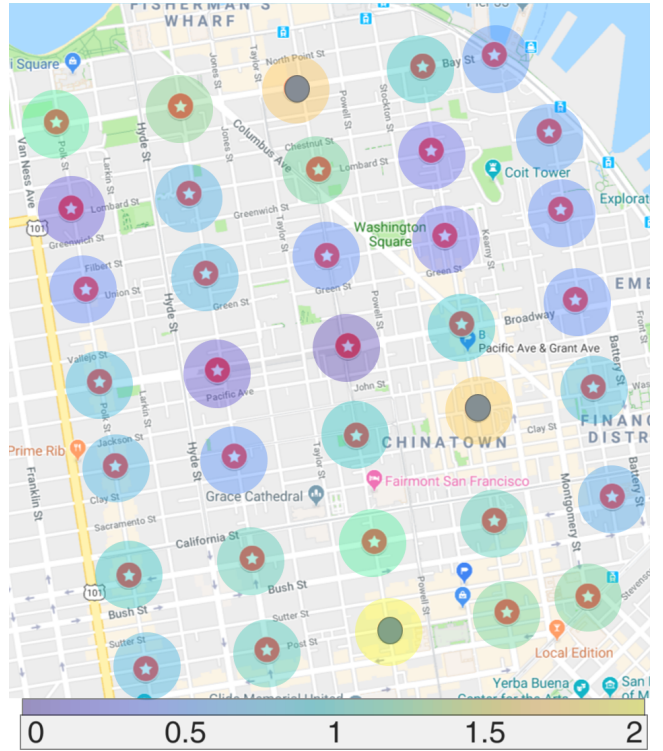


Fig. 3. The distribution of utility $U$.

$U(s)$, assigned to each location s. In Figure 3, $U(s)$ is shown with a logarithmic scale for better illustration. As the color bar at the bottom of the figure indicates, the closer the color at each location is to the right side of this bar, the higher patrol hours are assigned. For example, the state at $(3, 1)$ (the third state from the first row), where $C(s) = 48$ gets assigned patrol hour equals 106 which is approximately 2 in logarithmic scale. Therefore, its color is yellow in Figure 3 as indicated by the right tip of the color bar. Together with Figure 2, it can be observed that sensitive places and places with a higher number of crimes get assigned more patrol hours. Consequently, the rewards at those states are fairly low to discourage the criminal from visiting it. The cost at each location is proportional to the crime rate and inversely proportional to the police patrol hours. The patrol hours assigned to each place intends to minimize the expected cost incurred by the human adversary.

## VII. Conclusion

This paper introduces a general framework for deceiving adversaries with bounded rationality in terms of the obtained reward minimization. Leveraging the cognitive bias of the human from well-known prospect theory, we formulate the reward-based deception as a resource allocation problem in Markov decision process environment and solve as a signomial program to minimize the adversary's expected reward. We use police patrol hour assignment as the illustrative example and show the validity of our propose solution approach. It opens doors for further research on the topic to consider the scenarios where defender can move around and

react to the human adversaries in real time, and the human adversary has a learning capability to adapt the defender's deceiving policy.

## REFERENCES

[1] M. H. Almeshekah and E. H. Spafford, "Cyber security deception," in *Cyber deception*. Springer, 2016, pp. 23–50.

[2] J. Pawlick, E. Colbert, and Q. Zhu, "A game-theoretic taxonomy and survey of defensive deception for cybersecurity and privacy," *arXiv preprint arXiv:1712.05441*, 2017.

[3] S. Bonetti, "Experimental economics and deception," *Journal of Economic Psychology*, vol. 19, no. 3, pp. 377–395, 1998.

[4] T. Holt, *The deceivers: Allied military deception in the Second World War*. Simon and Schuster, 2010.

[5] E. Morgulev, O. H. Azar, R. Lidor, E. Sabag, and M. Bar-Eli, "Deception and decision making in professional basketball: Is it beneficial to flop?" *Journal of Economic Behavior & Organization*, vol. 102, pp. 108–118, 2014.

[6] M. L. Puterman, *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons, 2014.

[7] K. Horák, Q. Zhu, and B. Bošanský, "Manipulating adversary's belief: A dynamic game approach to deception by design for proactive network security," in *International Conference on Decision and Game Theory for Security*. Springer, 2017, pp. 273–294.

[8] H. A. Simon, "Models of man; social and rational." 1957.

[9] C. F. Camerer, *Behavioral game theory: Experiments in strategic interaction*. Princeton University Press, 2011.

[10] R. Yang, B. Ford, M. Tambe, and A. Lemieux, "Adaptive resource allocation for wildlife protection against illegal poachers," in *Proceedings of the 2014 international conference on Autonomous agents and multiagent systems*. International Foundation for Autonomous Agents and Multiagent Systems, 2014, pp. 453–460.

[11] C. Zhang, A. X. Jiang, M. B. Short, P. J. Brantingham, and M. Tambe, "Defending against opportunistic criminals: New game-theoretic frameworks and algorithms," in *International Conference on Decision and Game Theory for Security*. Springer, 2014, pp. 3–22.

[12] B. Wu and H. Lin, "Privacy verification and enforcement via belief abstraction," *IEEE Control Systems Letters*, vol. 2, no. 4, pp. 815–820, Oct 2018.

[13] P. Masters and S. Sardina, "Deceptive path-planning," in *Proceedings of the 26th International Joint Conference on Artificial Intelligence*. AAAI Press, 2017, pp. 4368–4375.

[14] D. Kahneman, *Thinking, Fast and Slow*. Macmillan, 2011.

[15] A. Tversky and D. Kahneman, "Advances in prospect theory: Cumulative representation of uncertainty," *Journal of Risk and uncertainty*, vol. 5, no. 4, pp. 297–323, 1992.

[16] D. Kahneman and A. Tversky, "Prospect theory: An analysis of decision under risk," in *Handbook of the fundamentals of financial decision making: Part I*. World Scientific, 2013, pp. 99–127.

[17] E. Norling, "Folk psychology for human modelling: Extending the bdi paradigm," in *Proceedings of the Third International Joint Conference on Autonomous Agents and Multiagent Systems-Volume 1*. IEEE Computer Society, 2004, pp. 202–209.

[18] P. B. Reverdy, V. Srivastava, and N. E. Leonard, "Modeling human decision making in generalized gaussian multiarmed bandits," *Proceedings of the IEEE*, vol. 102, no. 4, pp. 544–571, 2014.

[19] D. S. Hochbaum, "Complexity and algorithms for nonlinear optimization problems," *Annals of Operations Research*, vol. 153, no. 1, pp. 257–296, 2007.

[20] S. Boyd, S.-J. Kim, L. Vandenberghe, and A. Hassibi, "A tutorial on geometric programming," *Optimization and Engineering*, vol. 8, no. 1, 2007.

[21] B. An and M. Tambe, "Stackelberg security games (ssg) basics and application overview," *Improving Homeland Security Decisions*, p. 485, 2017.

[22] M. Jain, J. Tsai, J. Pita, C. Kiekintveld, S. Rathi, M. Tambe, and F. Ordónez, "Software assistants for randomized patrol planning for the lax airport police and the federal air marshal service," *Interfaces*, vol. 40, no. 4, pp. 267–290, 2010.

[23] W. B. Haskell, D. Kar, F. Fang, M. Tambe, S. Cheung, and E. Denicola, "Robust protection of fisheries with compass." 2014.

[24] R. Yang, C. Kiekintveld, F. OrdóñEz, M. Tambe, and R. John, "Improving resource allocation strategies against human adversaries in security games: An extended study," *Artificial Intelligence*, vol. 195, pp. 440–469, 2013.

[25] D. Kar, F. Fang, F. Delle Fave, N. Sintov, and M. Tambe, "A game of thrones: when human behavior models compete in repeated stackelberg security games," in *Proceedings of the 2015 International Conference on Autonomous Agents and Multiagent Systems*. International Foundation for Autonomous Agents and Multiagent Systems, 2015, pp. 1381–1390.

[26] S. Boyd and L. Vandenberghe, *Convex optimization*. Cambridge university press, 2004.

[27] Y. D. Abbasi, M. Short, A. Sinha, N. Sintov, C. Zhang, and M. Tambe, "Human adversaries in opportunistic crime security games: Evaluating competing bounded rationality models," in *Proceedings of the Third Annual Conference on Advances in Cognitive Systems ACS*, 2015, p. 2.

[28] K. Doya, "Modulators of decision making," *Nature neuroscience*, vol. 11, no. 4, p. 410, 2008.

[29] M. B. Short, M. R. D'orsogna, V. B. Pasour, G. E. Tita, P. J. Brantingham, A. L. Bertozzi, and L. B. Chayes, "A statistical model of criminal behavior," *Mathematical Models and Methods in Applied Sciences*, vol. 18, no. supp01, pp. 1249–1267, 2008.

[30] R. E. Moore, "Global optimization to prescribed accuracy," *Computers & mathematics with applications*, vol. 21, no. 6-7, pp. 25–39, 1991.

[31] E. L. Lawler and D. E. Wood, "Branch-and-bound methods: A survey," *Operations research*, vol. 14, no. 4, pp. 699–719, 1966.

[32] P. J. Brantingham and G. Tita, "Offender mobility and crime pattern formation from first principles," in *Artificial crime analysis systems: using computer simulations and geographic information systems*. IGI Global, 2008, pp. 193–208.

[33] E. Burnell and W. Hoburg, "Gpkit software for geometric programming," https://github.com/convexengineering/gpkit, 2018, version 0.7.0.

[34] M. ApS, *The MOSEK optimization toolbox for PYTHON. Version 7.1 (Revision 60)*, 2015. [Online]. Available: http://docs.mosek.com/7.1/quickstart/Using_MOSEK_from_Python.html