# Analyzing Human and Virtual Agent Interaction Under Irrational Decision Making

Mohammad H. Mamduhi* and Sandra Hirche*
*Department of Electrical Engineering and Information Technology
Technische Universität München, München 80333
Email: {mamduhi,hirche}@lsr.ei.tum.de

*Abstract*—The incapability of virtual agents to analyze the irrational behavior in human-virtual agent interactions urges to develop efficient mechanisms to capture the irrationality. Irrational actions in cognitive interactions are employed to attract the opponent's trust or to make future actions harder to anticipate. Recent works mostly investigate interactions of multiple decision makers with perfect knowledge which select optimal actions, a scenario that is not applicable when irrationality comes in. In this paper, the social interaction between a human and a virtual agent over a known utility function is considered while irrational actions are allowed. The existence of correlated equilibria is shown under some mild assumptions. Using the action history of both players, a measure of irrationality, which is variance-optimal, is proposed. Exploiting the irrationality measure, a probabilistic mechanism for decision making is suggested and it is shown that the sequence of decisions chosen under this mechanism leads to equilibria for sufficient number of iterations. Restating the interaction as an iterative game, the existence of equilibria is shown for two types of games i.e. simultaneous games and leader-follower (Stackelberg) games. Even though our method considers human-virtual agent interaction, the results extend to any two-player finite discrete game. An HCI experiment validates our proposed approach and corroborates the efficiency of the proposed method.

## I. INTRODUCTION

Technological advances in intelligent machines such as computers, robots, and virtual agents have led to an increased interest in forming joint tasks by human. In such a partnership various simplifications are imposed due to the highly complex aspects of human characteristics i.e. full information access or rational decision makers. In real interactions though, these ideal assumptions do not often hold. Irrationality of decision makers is one of the unexplored aspects of social interactions in majority of studies, and is rather untouched from the birth of the concept in half a century ago.

The early works utilizing virtual environments started in 60's by introducing graphical interfaces, operation systems, and physical devices. During recent years though, focus is turned to physical and cognitive interactions between human and intelligent machines [7] and [9]. Human-machine interaction is studied from different perspectives, such as social human behavior and decision making, establishment of different ways of communication and data interpretation, and implementing proper autonomous behaviors for physical and verbal/non-verbal interaction with the human. Numerous psychological and neurological works consider the first aspect and try to recognize human decision making mechanisms and map them out to intelligent machines, see e.g. [8] and [3]. In [4], studies show that human naturally tends to share the information resources and cooperate with the partners during an interaction, but not always doing that. Moreover, human mostly starts with a positive belief about the other agents, so tries to be cooperative unless receives any contentious reaction. Significant efforts are made on how to establish data communication and interpretation between human and autonomous agents. Audio-visual communication with the human is widely considered, e.g. in [10]. More recently, physical interaction between human and robots received a substantial interest [12], [6]. Recently, Game Theory is emerged as a powerful tool to model human-machine interaction, investigate the agents' decisions, their consequences on utility, and strategy management. In [2] a game theoretic solution in competitive social interaction is found by developing continuous versions of Prisoner's Dilemma[1] and Rope-Pulling[2] games. A machine-learning approach to identify human preferences considering social factors in one-shot games is proposed in [13]. In majority of studies about human-virtual agent interaction, rationality assumption holds to simplify the interpretation of human behavior, since irrationality makes it excessively hard to anticipate human intention during an interaction. In [14], a computational model of emotion-focused irrational behavior is developed to design virtual humans. Qualitative discussions in [15] investigate the role of irrational factors, but in strategic and financial management. Rationality assumption may seem reasonable in cooperation, but is restrictive in non-cooperative interactions. The necessity of considering irrationality is of great interest in economics, but also in any competitive game scenarios. The design and implementation of intelligent virtual agents capable of coping with irrationality is of high importance for advance human-computer interaction.

The contribution of this paper is to propose a methodology to establish the interaction under possible irrationalities, and study the existence of equilibria. We first remodel the interaction as a game, then prove the existence of equilibria for general interactions in presence of irrational actions under some mild assumptions. We introduce an approach for measuring irrationality, and exploit it to construct a Probability

---

[1]PD is a two-player non-zero-sum game revealing the importance of cooperation in iterative interactions. It shows that cooperation is more tied to mutual trust than to utility function.

[2]For more description see [1].

Distribution over the virtual agent's strategy set. Subsequently, we prove the sequence of decisions taken under the *pdf* is the optimal reaction against irrationality. Finally, an experiment in a human-virtual agent game corroborates the results.

The reminder of the paper is structured as follows: In section II, we state the problem, mathematical notation, and assumptions. Section III derives the game theoretical model of the interaction. Existence of equilibria for rational interaction is shown in section IV. The main contribution of the paper is presented in section V by introducing the irrationality measure and proposing a probability distribution under which the sequence of decisions made by computer is at equilibrium. Section VI validates the approach by an experiment.

## II. PROBLEM STATEMENT, NOTATION AND ASSUMPTIONS

Consider a discrete time finite human-virtual agent interaction over a known utility function . Each player competes to get the most out of the utility by selecting proper strategies from a finite strategy set $S$. At every iteration, the awarded utilities are mappings from both players' strategies into the set of positive real numbers, i.e. $J_C^i(S_C^i, S_H^i): S \times S \to \mathbf{R}^+$ and $J_H^i(S_C^i, S_H^i): S \times S \to \mathbf{R}^+$, where $S_C^i$ and $S_H^i$ denote the selected strategies at time step $i$. The objective is to construct a proper strategy profile for the virtual agent against the human, who also tries to maximize the own reward, such that within a finite number of iterations the maximum utility is achieved. Notably, both players being maximizers does not necessarily mean that their interaction is competitive. Since both players jointly determine the individual utilities, the maximum utility might be achieved by cooperation. Here rationality plays an important role to build up the reciprocal trust which is the essential factor for cooperation. If rationality is violated, the mutual trust between the players would easily be faded and cooperation stops. Rationality imposes two assumptions on strategic behaviors [5]: first, players make decisions based on some belief about opponent's behavior. Second, belief should be consistent with other players being rational, being aware of each other's rationality and so on, in an infinite cycle. Acting irrationally decreases the level of trust between the players and leads them to employ more competitive strategies.

The following assumptions enable a simplified decision making design to capture the irrationality by the virtual agent:
$A_1$: Players are able to keep the action history of each other during earlier iterations.
$A_2$: Strategy and action sets are finite and bounded.
$A_3$: Players are allowed either to employ a pure strategy from the strategy set $S$ or to mix between possible strategies, rationally or irrationally at each iteration.
$A_4$: Players are rational in employing irrational strategies. It means irrationality is purposely employed to mislead the opponent, and irrational actions are chosen according to some belief about the opponent's behavior. This assumption clarifies the difference between an *irrational action*, and a *foolish action*. In a foolish action the actor deviates from rationality having no specific reason behind.

The above assumptions are not severely restrictive. For example $A_2$ declares that the players have limited alternatives to choose and act based upon, which seems quite reasonable.

*Remark 1*: A direct consequence of utilities being known for both players at each iteration is that they are able to easily detect the irrational actions employed by each other.

## III. INTERACTION AS A GAME

Human-Computer Interaction, as a special form of human-virtual agent interaction, can be effectively restated within the game theoretic framework. Variety of games then exist theoretically correspond to various class of bilateral interactions. Here, we consider two types of dynamic interactions; one at which players act simultaneously, and the other with players acting sequentially. In what follows, the game theoretic models for both scenarios are derived, and well-defined games which reflect the specific properties of each type are introduced. First, we start with some basic game theoretic definitions [1].

*Definition 3.1: Zero-Sum vs. Non-Zero-Sum Games*: In *zero-sum* games, sum of the players' gain and loss is zero, while in *non-zero-sum* games this sum is a nonzero constant. Notably, cooperation happens in non-zero-sum games only, since there is no obligation on players' costs and rewards to be identical. In zero-sum games only competition occurs.

*Definition 3.2: Finite vs. Infinite Games*: In a *finite* game, each player selects a strategy from a limited number of alternatives, while in *infinite* games, infinite choices are possible to be selected.

*Definition 3.3: Simultaneous vs. Leader-Follower*: In *simultaneous* games, players make their decisions simultaneously at each iteration such that neither side could would know it as a priori. In contrast, the game in which one player, called *leader*, declares his strategy first and the other player, called *follower*, has to react to the leader's strategy, is called *leader-follower*, or *Stackelberg* game.

In the following sections, we focus on non-zero-sum games which enable us to consider cooperation in the interactive loop. Moreover, only finite games are considered according to assumption $A_2$, which automatically excludes the continuous-kernel games from the discussions. In the next section, we preliminary study the existence of equilibria in constant-sum finite simultaneous games, and then continue with its Stackelberg counterpart both under rationality assumption.

## IV. NASH EQUILIBRIUM IN RATIONAL INTERACTION

Before going through the irrationality in the games of interest, we shortly review the existence of equilibria in rational interactions. In the following, the existence of Nash equilibrium for two-player non-zero-sum finite simultaneous and Stackelberg games with rational players are summarized. By relaxing rationality though, the results are no longer valid.

## A. Equilibrium Strategies in Finite Simultaneous Games

To make the discussions in this section explicit, the difference between two mathematical representations of finite games is clarified first, i.e. *normal* and *extensive* representations. An extensive description explicitly reflects the correspondences between different strategies and their outcomes, and provides primary information such as order of play, dynamics evolution, or information exchange. A normal description however, is a compact form which suppresses the dynamic character of a game. Conclusively, a game with same strategies, outcomes, and equilibria can be represented by either model.

*Definition 4.1: Non-Cooperative Equilibrium*: In an iterated two-player non-zero-sum finite game, the sequence of discrete strategies $\{S_H^{1^*}, \ldots, S_H^{K^*}, S_C^{1^*}, \ldots, S_C^{K^*}\}$, with final iteration $K$, is a non-cooperative equilibrium if the following two inequalities are satisfied for all the available strategies:

$$J_H^* = J_H(S_H^*, S_C^*) \geq J_H(S_H, S_C^*)$$
$$J_C^* = J_C(S_H^*, S_C^*) \geq J_C(S_H^*, S_C) \tag{1}$$

where, $S^* = \{S^{1^*}, \ldots, S^{K^*}\}$ is the profile of dominating strategies. $K = 1$ defines the static equilibrium strategy [1].

*Proposition 1*: Every $N$-player non-zero-sum finite sequential game in extensive form admits a non-cooperative equilibrium solution in mixed strategies. (Proof in [1]).

*Proposition 2*: Every $N$-player finite sequential game in normal form necessarily admits a *perfect equilibrium*, and every perfect equilibrium is a Nash equilibrium. (Proof in [1]).

Therefore, competitive finite simultaneous games with rational players necessarily admit a Nash equilibrium, either pure or mixed. Furthermore, rational cooperation always improves the outcome of a game compared to the Nash equilibrium [11]. In other words, under the worst case possible situation that players make no cooperation, Nash equilibrium still exists.

## B. Stackelberg Strategies in Leader-Follower Games

We first introduce the concept of Stackelberg equilibrium, and then discuss the existence of equilibria in finite Stackelberg games. Human ($H$) and computer ($C$) represent the leader and follower, respectively.

*Definition 4.2*: Let $S$ denotes the pure strategy space of the players in a finite Stackelberg game, and $J_H(S_H, S_C)$ and $J_C(S_H, S_C)$ represent the utilities incurred to them, corresponding to strategy pair $\{S_H, S_C \in S\}$. Then set $R^C(S_H) \subset S$ defined for each $S_H \subset S$ by

$$R^C(S_H) = \{\bar{S}_C \in S : J_C(S_H, \bar{S}_C) \geq J_C(S_H, S_C), \ \forall S_C \in S\}$$

is the optimal response (rational reaction) set of the follower to the strategy $S_H \in S$ of the leader.

*Definition 4.3*: In a two-player finite Stackelberg game, strategy $S_H^* \in S$ is called a Stackelberg equilibrium strategy for the leader, if

$$J_H(S_H^*, S_C^*) = \max_{S_H} J_H(S_H, S_C^*) \tag{2}$$

*Proposition 3*: Every two-player finite Stackelberg game admits a Stackelberg equilibrium for the leader.

*Proof:* : Since $S$ is a finite set, and $R^C(S_H)$ is a subset of $S$ for each $S_H \in S$, the proof readily follows from (2). ∎

*Definition 4.4*: Let $S_H^*$ be a Stackelberg strategy for the leader. Then any element $S_C^* \in R^C(S_H^*)$ is an optimal strategy for the follower that is in equilibrium with $S_H^*$. The pair $\{S_H^*, S_C^*\}$ is a Stackelberg solution of the game, and cost pair $(J_H(S_H^*, S_C^*), J_C(S_H^*, S_C^*))$ is the corresponding Stackelberg outcome.

*Remark 2*: Uniqueness of equilibria in either simultaneous or Stackelberg game is not of interest, since in the latter case, the leader's utility is unique even if the strategy is not. For simultaneous games also, uniqueness can be discussed only in very simple class of information patterns. So, it is generally infeasible to be studied, especially when decisions are made based on a probabilistic chance mechanism.

## V. A MEASURE OF IRRATIONALITY AND EXISTENCE OF CORRELATED EQUILIBRIA

This section offers the main contribution of the paper by studying the reactive policy of computer against the irrationality in human behavior. At first, a probabilistic measure is proposed to quantify the level of rationality in human behavior. Then, we show if computer react against the human actions by mixing between the available strategies, according to a well-defined *pdf*, the correlated equilibria exist.

## A. A Measure of Irrationality

It is essential for computer to have a quantitative understanding of human rationality. To satisfy this, a cumulative variance-optimal probabilistic measure is proposed based on the mutual action history of the players. Since, the human decision making mechanism depends on a lot of internal factors, and cannot be deterministically captured by mathematical models, we model the human policies by discrete random variables chosen under any proper probabilistic mechanism. Computer, as human opponent, needs to successfully measure how far the actual behavior of the human is from the rational behavior. Then, computer employs proper strategies to react against human behavior in future. The following theorem facilitates this by introducing a method to optimally measuring the distance of two random variables.

*Theorem 5.1*: Suppose that computer's expectation of human rational strategy at $i$'th iteration of the game, given the utility at steak, is defined by $E^C[\hat{S}_H^i | J_H^i]$. Let $S_H^i$ represent the actual strategy taken by human, either rational or irrational. Then, the smallest possible mean square error resulting from computer's expectation of human's rational strategy and human's actual strategy at time step $i$ is

$$E[(S_H^i - E^C[\hat{S}_H^i | J_H^i])^2] \tag{3}$$

where, $J_H^i$ is utility function for human at $i$'th iteration.

*Proof:* We show that the square error in (3) is the smallest quadratic error function between human's actual strategy and computer's expectation of human's rational strategy than any other function $g$ expecting the latter. As already stated, $S_H^i$ and $J_H^i$ are discrete random variables. Assume their Probability Mass Functions are represented by $P_{S_H}(S_H^i)$ and $P_{J_H}(J_H^i)$, with conditional PMF $P_{S_H|J_H}(S_H^i|J_H^i)$. Suppose $g : \mathbf{R} \to \mathbf{R}$ is any measurable function i.e. $S_H^i g(J_H^i)$ is either nonnegative or integrable. Then, for any specific event $J_H = J_H^i$ such that $P_{J_H}(J_H^i) > 0$ we have

$$
\begin{aligned}
&E[E[S_H|J_H]g(J_H)] \\
&= \sum_{J_H^j} E[S_H|J_H = J_H^j]g(J_H^j)P_{J_H}(J_H^j) \\
&= \sum_{J_H^j}\sum_{S_H^j} S_H^j P_{S_H|J_H}(S_H^j|J_H^j)g(J_H^j)P_{J_H}(J_H^j) \\
&= \sum_{S_H^j, J_H^j} S_H^j g(J_H^j)P_{S_H, J_H}(S_H^j, J_H^j) \\
&= E[S_H g(J_H)].
\end{aligned}
$$

Now, suppose $E[(S_H^j)^2] < \infty$. Since, $E[S_H|J_H]$ is a function of $J_H$, $E[E[S_H|J_H]E[S_H|J_H]] = E[S_H E[S_H|J_H]]$, then:

$$
\begin{aligned}
&E[(S_H - g(J_H))^2] \\
&= E[(S_H - E[S_H|J_H])^2] \\
&+ E[(E[S_H|J_H] - g(J_H))^2] \\
&+ E[(S_H - E[S_H|J_H])(E[S_H|J_H] - g(J_H))] \\
&\geq E[(S_H - E[S_H|J_H])^2],
\end{aligned}
$$

where the inequality is ensured considering the facts that $E[(E[S_H|J_H] - g(J_H))^2] \geq 0$ and the term $E[(S_H - E[S_H|J_H])(E[S_H|J_H] - g(J_H))]$ is zero. ∎

Exploiting theorem 5.1, we introduce the following variance-optimal measure which enables computer to quantify human's rationality during the interaction:

$$
\sigma_{\hat{S}_H^i}^C = \frac{1}{n}\sum_{j=1}^n E[(S_H^j - E^C[\hat{S}_H^j|J_H^j])^2]. \tag{4}
$$

(4) measures the average mean square error between human's actual and rational strategies from the initial time step up to the current $i$'th one. The measure is applicable for both games considered in this paper, with only difference being $n = i - 1$ in simultaneous games, and $n = i$ in Stackelberg games.

*Theorem 5.2*: For sufficiently large number of interactive iterations in a finite two-player game, the irrationality measure $\sigma_{\hat{S}_H^i}^C$ starts permanently decreasing after a finite iterations.

*Proof:* Consider the situation under which the human selects his strategies rationally at all the iterations, then $E[S_H^j] = E^C[\hat{S}_H^j|J_H^j]$ and according to (4), $\sigma_{\hat{S}_H^i}^C$ is zero. On the other hand, every irrationally selected strategy leads to quadratic expected error $E[(S_H^j - E^C[\hat{S}_H^j|J_H^j])^2] > 0$ at corresponding time step $j$. Since, both $S_H$ and $J_H$ are bounded finite discrete

random variables, then $E[(S_H^j - E^C[\hat{S}_H^j|J_H^j])^2] < \infty$ and eventually there exists a finite integer $m > 0$ such that

$$
E[(S_H^k - E^C[\hat{S}_H^k|J_H^k])^2] \leq \sum_{j=1}^{k-1} E[(S_H^j - E^C[\hat{S}_H^j|J_H^j])^2]
$$

for all $k \geq m$, implying that $\sigma_{\hat{S}_H^i}^C$ is a decreasing sequence for all $k \geq m$. Moreover, $\sigma_{\hat{S}_H^i}^C$ is a positive sequence, then

$$
\lim_{n \to \infty} \sigma_{\hat{S}_H^i}^C \to 0
$$

which declares that, if a finite game with possibility of irrational actions under $A_1$-$A_4$ assumptions repeats infinitely, the error between computer's expectation of human's rational and actual strategy tends to zero. It means, although human can employ irrational strategies, computer can perfectly predict it. If $\sigma_{\hat{S}_H^i}^C \to 0$ applies, the equilibria studied in Section IV are again feasible. ∎

Qualitatively explaining, by every irrational action human shows, the computer is informed about human decision making mechanism. In other words, their information sets having no intersection initially, start to possess some common space with each non-zero $\sigma_{\hat{S}_H^i}^C$. Then, after a finite number of iterations ($m$), the inaccuracy of computer in expecting human's next move will be less than all the previous ones, which have been kept in computer's memory. Clearly, this specific iteration would be further if human employs less irrational strategies.

*Remark 3*: The necessity of the assumption $A_2$ can be seen from the theorem 5.2. Otherwise, the monotonicity of $\sigma_{\hat{S}_H^i}^C$ cannot be ensured with finite iterations.

### B. Existence of Correlated Equilibria

Here we show the existence of *Correlated Equilibria* for two-player non-zero-sum finite simultaneous and Stackelberg games with possibility of irrational actions and under the assumptions $A_1 - A_4$. We derive the results only for simultaneous games, since showing the existence of equilibria for Stackelberg games is exactly the same. The formal definition of correlated equilibria is as follows [5]:

*Definition 5.1: Correlated Equilibria*: Consider a finite game $G$ with $J_i(s_i, s_{-i})$ denotes the utility function for player $i$, where $s_{-i} = [s_j]_{j \neq i}$ represents the strategy profile of all players except $i$'th one. Suppose $\Delta(S)$ is the set of all discrete probability measures over the set $S$, then the joint distribution $\pi \in \Delta(S)$ is a correlated equilibrium if and only if

$$
\sum_{s_{-i}} \pi(s)[J_i(s_i, s_{-i}) - J_i(\bar{s}_i, s_{-i})] \geq 0 \tag{5}
$$

for all players $i$ and all $s_i, \bar{s}_i \in S$ i.e. $s_i \neq \bar{s}_i$.

Next theorem shows $\pi(s)$ is a correlated equilibrium if it represents a discrete normal distribution over the joint strategy set $S$ with standard deviation $\sigma_{\hat{S}_H^i}^C$ and mean $\mu^i$ defined in (6).

*Theorem 5.3*: In a two-player finite simultaneous game, with possibility of irrational actions from human, $\pi^*(s) :=$

$\mathcal{N}(\mu^i, \sigma^C_{\hat{S}^i_H})$ distributed over the strategy set $S$ is a correlated equilibrium in response to human strategies, with $\mu^i$ defined as

$$\mu^i = \omega^i_1 S_{C_1} + \omega^i_2 S_{C_2} + \ldots + \omega^i_L S_{C_L} \qquad (6)$$

where, $\omega^i$s are weights assigned to each strategy $S_{C_l} \in S$ according to priority of that strategy for computer, and $\omega^i_1 + \ldots + \omega^i_L = 1$ at each time step $i$, with $L$ being the finite number of available strategies.

*Proof:* Let $\Pi(s)$ be the set of all normal distributions over the set $S$, with different standard deviations, including $\pi^*(s)$. Since $\omega_l$s can be regulated by computer independent of human at each time step, according to the more preferred strategy to respond human decisions, let assume all the normal distributions in $\Pi$ have the same mean for each specific strategy profile of the human. Now, we need to show that the following inequality holds for all the strategies $S^i_H$:

$$\sum_{S^i_H} \pi^*(s)[J^i_C(\bar{S}^i_C, S^i_H) - J^i_C(S^i_C, S^i_H)] \geq 0$$

where, $\bar{S}^i_C \in \pi^*(s)$ is the strategy chosen under $\pi^*(s)$, $S^i_C \in \Pi(s) - \pi^*(s)$ is any strategy chosen under the other distributions in $\Pi(s)$ than $\pi^*(s)$. Now consider for a specific human strategy $S^i_H$, either rational or not, $J^i_C(\bar{S}^i_C, S^i_H) - J^i_C(S^i_C, S^i_H) < 0$. It means there should be a normal distribution $\bar{\pi}(s) \in \Pi$ by which computer gets more utility choosing its strategy profile under. Since, the mean values for all the normal distributions in $\Pi$ are identical for each specific $S^i_H$, and also, computer acts rationally, $\bar{\pi}(s)$ should have a better expectation of human's strategy. In the other words, $\bar{\pi}(s)$ should make a narrower distribution over the strategy set $S_C$ for computer, which practically means it should have smaller variance than $\pi^*(s)$. However, this is in contradiction with the theorem 5.1, since $\sigma^C_{\hat{S}^i_H}$ captures the smallest error of expecting human's strategy, therefore, from the moment that the sequence of $\sigma^C_{\hat{S}^i_H}$ starts permanently decreasing, which is a finite moment $m$ according to the theorem 5.2, $J^i_C(\bar{S}^i_C, S^i_H) - J^i_C(S^i_C, S^i_H) \geq 0$ always holds. On the other hand, $\pi^*(s)$ is a normal distribution which always assigns a non-negative probability to different strategies, so $\pi^*(s)[J^i_C(\bar{S}^i_C, S^i_H) - J^i_C(S^i_C, S^i_H)]$ is always non-negative for time steps greater than $m$, and hence (5) is fulfilled, and $\pi^*(s)$ is a correlated equilibrium. ∎

*Remark 4*: We do not claim $\pi^*(s)$ is the only equilibrium, or even it is the ever-best solution. We show a specific way to deal with irrationality in finite games, and prove the existence of equilibria by establishing a probabilistic mechanism.

*Remark 5*: The correlated equilibrium is a more general concept than Nash, i.e. every Nash equilibrium is a correlated equilibrium, but not the other way around. So, the type of equilibria we derived is weaker than the concept of Nash equilibrium. However, in these types of problems where cooperation exists, and rationality assumption is relaxed, looking for Nash equilibrium is impracticable.

The proof of theorem 5.3 can be repeated for the Stackelberg games, almost the same, with the only difference again the
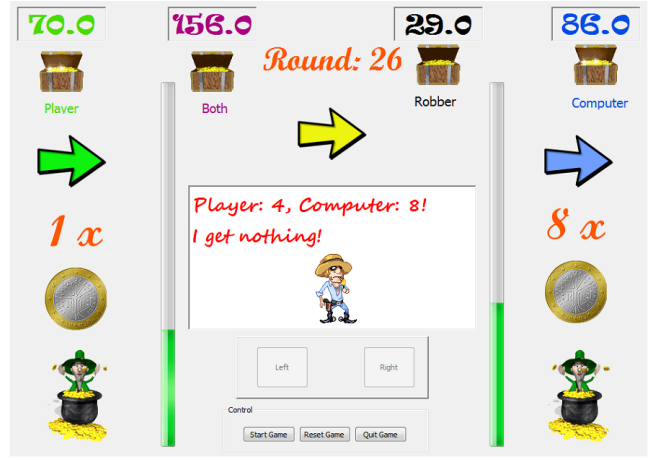


Fig. 1. Stackelberg Game Interface

upper limit of the summation in (4) i.e. $n = i$. In the next section, the results of an experiment are presented as a finite Stackelberg game, to validate the theorems.

## VI. Experiments and Results

In this section, we will validate our results in a two-player discrete finite binary decision making Stackelberg game, as depicted in Fig. 1. The game are played by a human as leader, and computer as follower. At the start of each new iteration, two coin values appear at either sides, and game's goal would be to collect as much as possible coins, based on the announced coin values. Every player has two alternatives to act and should choose between either *left* or *right*. Only if both of the players agree on one specific side, the utility on the corresponding chosen side will be awarded to them. If they do not agree on a particular side, then neither side's utilities would be awarded, therefore, the players need to cooperate to achieve some utilities. The utilities are awarded according to a rule, i.e. the whole amount of one side's utility is awarded to the player on the corresponding side, and the player on the other side gets half of this amount. For example, let human plays in the left side which value 10$ is appeared, and computer is in the right side with utility 4$. Thus, the rational action for both players is to cooperate and choose *left* side, so human and computer are awarded 10$ and 5$, respectively, while agreement on the other side would end them up with 4$ and 2$.

There are six different strategies ($l = 6$) under which the players can choose how to react with respect to each other, namely *Cooperative*, *Non-cooperative*, *Tit-for-Tat without Forgiveness*[3], *Tit-for-Tat with one Shot Forgiveness* (TFTF)[4], *Tit-for-Tat with a fixed Probability of Forgiveness*[5], and the last one is our derived strategy based on the irrationality

---

[3]Tit for Tat is a strategy under which a player starts with cooperation and will always cooperate unless provoked by the other player. If provoked the player will retaliate.

[4]TFTF considers forgiveness just once in response to non-cooperative action of the opponent.

[5]This strategy considers a fixed probability to forgive non-cooperative behavior of the opponent at each iteration of the game.

## TABLE I
### HCI in a Binary-Decision-Making Game and Utility Awarded Under Different Strategies

| Strategies | Irrational Actions | $\sum_{i=1}^{30} J_H^i$ | $\sum_{i=1}^{30} J_C^i$ |
|---|---|---|---|
| Cooperative | 23.3% | 71.7% | 58% |
| Non-cooperative | 16.7% | 57.7% | 46% |
| Tit-for-Tat | 20% | 60.8% | 59.4% |
| TFTF | 20% | 60% | 58% |
| 20% TFTF | 26.7% | 53.5% | 56% |
| $\pi^*(s)$ | 20% | 61.5% | 62.4% |
| Rational Nash Eq. | 0% | 69.6% | 67.3% |

measure in (4). To compare the equilibria obtained under different strategies with the pure Nash equilibrium, the latter is also calculated. Clearly, N.E. is superior since it takes place under rationality assumption. In table I, we demonstrate that interacting human with computer under the proposed decision making mechanism $\pi^*(s)$, and with rationality assumption relaxed, provides better solution than employing pure strategies. We have repeated each game for 30 iterations, and utilities are bounded between $[0 - 8]$\$ and randomly distributed for different time steps. Moreover, each game is played by ten people, each five times. The participants were given the basic information about the game and that they can employ irrational decision wherever they could take advantage of, but they did not know about the strategies taken by the computer.

As it can be seen, computer's final utility is maximized by choosing the strategies, against human's irrationality, under the *pdf* $\mathcal{N}(\mu^i, \sigma_{\hat{S}_H^i}^C)$ distributed over the whole set of feasible strategies, and with all the different strategies weighted uniformly. The "Irrational Actions" column shows how much human acts irrationally at those 30 iterations. $\sum_{i=1}^{30} J_C^i$ and $\sum_{i=1}^{30} J_H^i$ are also the achieved values out of the maximum achievable utility for each player. According to the first row in Table I, when computer acts based on cooperation, the most difference between their utilities occurs. It is sometimes called *Blind Cooperation*. TFT strategies however, are more reasonable, since it reacts to human's unfriendly actions. This is shown by making the achieved utilities closer to each other.

Here the results are a small step towards studying irrationality in more complex interactions between human and intelligent machines. Human-involved interactions without considering the possibility of strange behaviors from human is far from reality, and necessity of more studies to detect and manage those behaviors, specially when intelligent machines try to follow the human decision making mechanism, seems inevitable. We show that for simple interactions, irrationality can be treated, however, dealing with more complicated interactions such as human-robot interaction, and continuous interactions need much more studies.

## VII. Conclusion and Future Works

In this paper, we investigate the human-virtual agent interaction by taking the irrationality in decision making into account. We first remodel the interaction as a two-player game, and discussed the existence of equilibria holding the rationality assumption. Then, a variance-optimal approach to measure the irrationality of human's behavior is proposed which is based on the interaction's history and known utility function. Afterwards, we show that this measure is a decreasing sequence for sufficiently large number of interactive iterations. Finally, the existence of correlated equilibria is proved by introducing a probabilistic mechanism, and experiments validate results.

Since, irrationality in HCI and more generally in game theory has not been considerably discussed, it could potentially be a novel research line. As future works, the extension from discrete decision making process to continuous interaction, including multi-modal communication is planned. Developing new experimental setup, such that the game can be played by haptic devices, will also be a further step helping a lot in studying continuous interactions.

## References

[1] T. Basar and G. J. Olsder, "Dynamic Noncooperative Game Theory", 2nd Ed.,*Academic Press*, 1998, CA, USA, pp. 3, 4, 10-11, 99, 127-129.

[2] D. A. Braun, P. A. Ortega, D. M. Wolpert, "Nash Equilibria in Multi-Agent Motor Interactions", *PLoS Computational Biology*, August 2009, Vol. 5, Issue 8, e1000468.

[3] M. Cao, A. Stewart, N. E. Leonard, "Convergence in Human Decision-Making Dynamics", *System and Control Letters*, 2010, Vol. 59, pp. 87-97.

[4] P. F. Dominey, F. Warneken, "The basis of shared intentions in human and robot cognition", *New Ideas in Psychology*, 2009.

[5] D. Fudenberg, J. Tirole, "Game Theory", 7th Printing,*The MIT Press*, 2000, USA, pp. 48-49, 53-59.

[6] J. R. M. Hernandez, M. Lawitzky, A. Mortl, Dongheui Lee, and S. Hirche, "An Experience-Driven Robotic Assistant Acquiring Human Knowledge to Improve Haptic Cooperation", *IEEE International Conference on Intelligent Robots and Systems (IROS)*, 2011, pp. 2416-2422.

[7] A. Kucukyilmaz, S. O. Oguz, T. M. Sezgin, and C. Basdogan, "Improving Human-Computer Cooperation Through Haptic Role Exchange and Negotiation", Book Chapter, *To appear in A. Peer, and C. Giachritsis, Immersive Multimodal Interactive Presence (Ch. 13)*, Springer-Verlag.

[8] D. Lee, "Game theory and neural basis of social decision making", *Nature neuroscience*, April 2008, Vol. 11.

[9] N. Moray, "Identifying Mental Models of Complex Human-Machine Systems", *International Journal of Industrial Ergonomics*, 1998, Vol. 22, pp. 293-297.

[10] K. Nickel, H. K. Ekenel, M. Voit, and R. Stiefelhagen, "Audio-Visual Perception of Humans for a Humanoid Robot", *2nd International Workshop on Human-Centered Robotics Systems*, 2006, Munich, Germany.

[11] R. P. Nielsen, "Cooperative Strategy", *Strategic Management Journal*, 1988, Vol. 9, Issue 5, pp. 475-492.

[12] S. O. Ouguz,T. M. Sezgin, and C. Basdogan, "Supporting Negotiation Behavior in Haptics-Enabled Human-computer Interfaces", *Transactions on Haptics, Special Issue on Haptic Human Robot Interaction*, 2012.

[13] G. Yaakov, A. Pfeffer, F. Marzo, and B. J. Grosz, "Learning Social Preferences in Games", *In Proceedings, Nineteenth National Conference on Artificial Intelligence*, July 2004, California, USA.

[14] S. Marsella, and J. Gratch, "A step toward irrationality: using emotion to change belief", *Proceedings of the first international joint conference on Autonomous agents and multi-agent systems: part 1*, 2002, NY, USA.

[15] Y. Guo, "Study on Irrational Strategic Management", *In International Journal of Marketing Studies*, November 2009, Vol. 1, No. 2, China.