



Queen's Economics Department Working Paper No. 1033

## Rationalizing Irrational Beliefs

Geoffrey Dunbar  
Simon Fraser University

Juan Tu  
Queen's University

Ruqu Wang  
Queen's University

Xiaoting Wang  
Brock University

Department of Economics  
Queen's University  
94 University Avenue  
Kingston, Ontario, Canada  
K7L 3N6

2-2006

# Rationalizing Irrational Beliefs\*

Geoffrey Dunbar, Juan Tu, Ruqu Wang, and Xiaoting Wang

February 2006

Abstract. In this paper, we re-examine various previous experimental studies of the Centipede Game in the literature. These experiments found that players rarely follow the subgame-perfect equilibrium strategies of the game, and various modifications to the game were proposed to explain the outcomes of the experiments. We here offer yet another modification. Players have a choice of whether or not to believe that their opponents use subgame-perfect equilibrium strategies. We define a ‘behavioral equilibrium’ for this game. This equilibrium concept can reproduce the outcomes of those experiments.

*JEL Classification Numbers: C72, C91*

*Keywords: centipede game, game theory, experimental economics, behavioral economics*

---

\* Corresponding author: Ruqu Wang, Economics Department, Queen’s University, Kingston, Ontario K7L 3N6, Canada. Geoffrey Dunbar: Simon Fraser University; Juan Tu: Queen’s University; Xiaoting Wang: Brock University. Ruqu Wang’s research is supported by the Social Sciences and Humanities Research Council of Canada. We thank Jim Bergin and Lester Kwong for helpful comments.

# 1 Introduction

In dynamic games of perfect information, the concept of subgame-perfect equilibrium is most commonly used in the prediction of players' behavior. Consider a generic game of finitely many moves, the subgame-perfect equilibrium always exists and is unique. While the equilibrium concept is easily understood and the equilibrium characterization is usually straight-forward, challenges to its ability to predict players' behavior grow in the literature, both on theoretical front and experimental front.

Back in the nineteen eighties, Rosenthal (1982) constructed a game (which was later dubbed the "Centipede Game") which consisted a sequence of a hundred moves. In this game, each player moves at every alternative period, either to pass to the next period or to end the game right away. Passing the game to the next period would yield a larger total pile of money, but it strictly reduces the money one gets if the opponent ends the game then. The unique subgame perfect equilibrium (SPE) is that the first player ends the game at the first node and each player gets a small sum. Rosenthal argued that it is highly unlikely that, in practice, players will actually choose the SPE strategies when they play that game.

Since then, various centipede game experiments have been conducted to test the predictive power of the concept of SPE. McKelvey and Palfrey (1992) reported that only 15% of the players chose to end the game at the first node in the high-payoff version, and as little as 0.7% in other versions of the centipede game. In a much simplified two-move extensive form game, Goeree and Holt (2001) documented that players usually do not trust their rivals to be rational, and as a result, credible threats may not be credible at all.

In an attempt to reconcile the differences between the theory and the experiment outcomes, various modifications to the assumptions of the games used in the experiments have been proposed. McKelvey and Palfrey (1992), for example, propose that a player believes that the opponent is an altruist with some positive probability. They find that even a very small such probability can induce players to adopt mixed strategies in the early rounds of the game, mimicking the observed behaviors in their experiment. A few years later, McKelvey and Palfrey (1998) use a quantal choice model to re-examine the same experimental results. They show that if we assume that the probability of implementing a particular strategy is increasing in the equilibrium payoff of the strategy, then the observed behavior more or less coincides with the predictive behavior. Zauner (1999), on the other hand, adds yet another alternative to the explanation of McKelvey and Palfrey's experimental results by injecting a random perturbation to each player's payoffs. Various types of perturbations are explored and two best-fit models are selected.

At the same time, many game theorists have proposed alternatives to the basic as-

sumptions that lead to SPE, including common knowledge of rationality and backward induction. For example, Aumann (1992) formalizes the idea of higher order mutual knowledge. Caplan (2001) treats irrationality as a standard good, and players need to pay to get closer to some (irrational) “bliss belief”. Basu (1988) argues that each history of moves reveals certain characteristics of players to one another, and therefore the outcomes of a game depend on these revealed characteristics (instead of depending on rationality alone).

Recent advances in psychology have also helped in explaining why players in experiments behave differently from what SPE predicts. Epstein et al (1992) conduct studies to test the cognitive-experiential self-theory. They confirm that two conceptual systems, an experiential system and a rational system, each operate by its own rules of inference inside the same individual. To some extents, an individual can switch from one system to another. Tirole (2002) builds on similar psychological findings and explores their implications in an individual’s decision making process. He proposes a model of rational irrationality which can explain why people rehearse good news and selectively forget bad news, which is a universal behavior.

In this paper, we argue along the lines of the above psychological findings and propose yet another explanation on the “irrational behaviors” on the theoretical front. We emphasize on the observation that that even if all players understand fully the concept of subgame-perfect equilibrium and even if no players believe that other players are altruists, they still do not follow the SPE strategies when playing the centipede game. We assume that a player can choose to be “rational” or “behavioral”. If being “behavioral” yields better outcome than being “rational”, then a player would choose to “behavioral” (or, in terms of standard game theory terminology, “irrational”.) The intuition for this to happen is as follows. SPE strategies are optimal for a player only when other players follow them. If players do not believe that other players will follow SPE strategies, then those SPE strategies are not optimal anymore. In the model, we specify an alternative belief for each player regarding the behavior of other players. Each player then has a choice of selecting his belief (between the SPE strategy and the alternative one) at the beginning of the game and then optimizing given the selected belief. A “behavioral equilibrium” is formed if each player is better off in the actual outcomes by selecting the alternative belief. These outcomes of the game are determined by the strategies the players *actually* used in the game.

The basic idea behind the “behavioral equilibrium” concept is that players can choose to believe that their counterparts can be either fully rational (such that SPE strategies are the best response) or somewhat irrational (so that SPE strategies are not best response anymore). Given any belief, the players still optimize by choosing the best strategy. This is the same as in a subgame-perfect equilibrium. However, the difference between a behavioral equilibrium and a subgame-perfect equilibrium is that those alternative

beliefs in a behavioral equilibrium do not usually coincide with those players' actual strategies. If the two are the same, a subgame-perfect equilibrium is formed. Therefore, these alternative beliefs are somewhat irrational. Still, these irrational beliefs generate better payoffs than those SPE beliefs. Thus, players will choose these irrational beliefs rationally.

The rest of this paper is organized as follows. In Section 2, we analyze a few centipede games using the concept of "behavioral equilibria". In Section 3, we analyze some of the experiments in centipede games in the literature. In Section 4, we conclude.

## 2 Centipede Games and Behavioral Equilibria

We begin with a general description of the centipede games.

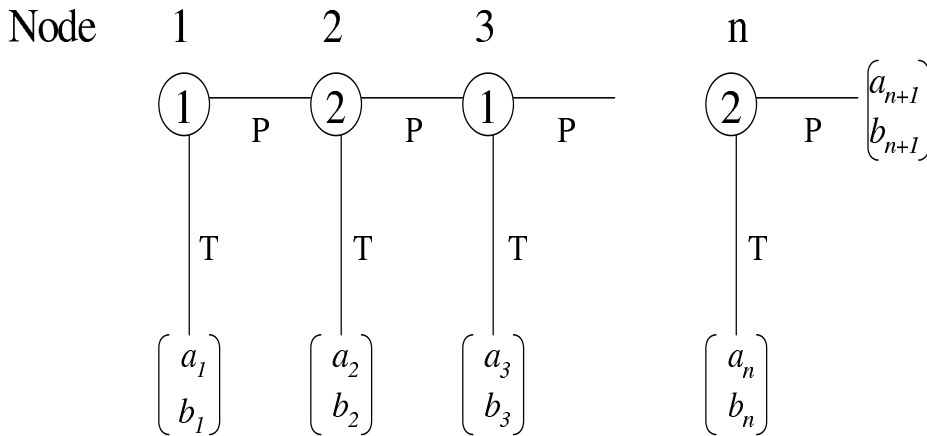


Figure 1: A General  $n$ -Move Centipede Game

There are two players, 1 and 2, playing the centipede game of  $n$  moves in Figure 1. To simplify notation, we assume that  $n$  is even.

In this game,  $a_1 > a_2, a_3 > a_4, \dots, a_{2i-1} > a_{2i}, \dots, a_{n-1} > a_n$ , and  $b_2 > b_3, b_4 > b_5, \dots, b_{2j-2} > b_{2j-1}, \dots, b_n > b_{n+1}$ . It is straight-forward to check that the unique subgame-perfect equilibrium strategy for each player is to play  $T$  whenever it is his turn to move. Given this strategy, the equilibrium outcome of the game is that player 1 plays  $T$  at the very beginning and ends the game with payoffs  $(a_1, b_1)$ .

Now suppose that before the start of the game, the two players choose a belief simultaneously. Player 1 secretly chooses a belief from  $\{SPE_1, B_1\}$ ; at the same time, player 2 secretly chooses a belief from  $\{SPE_2, B_2\}$ . Here,  $SPE_i$  represents player  $i$ 's subgame-perfect equilibrium belief on his opponent  $j$ 's behavior; i.e., player  $j$  will play  $T$  whenever it is his move. On the other hand,  $B_i$  denotes player  $i$ 's alternative belief. Let  $B_1 = (p_2, p_4, \dots, p_n)$  be player 1's belief, where  $p_{2k}$  is the probability that player 2 will play  $T$  at node  $2k$  conditional on node  $2k$  being reached. For SPE belief,  $SPE_1 = (1, 1, \dots, 1)$ . Similarly, we define  $B_2 = (p_1, p_3, \dots, p_{n-1})$ , and  $SPE_2 = (1, 1, \dots, 1)$ .

The subgame-perfect equilibrium belief  $SPE_i$  is the only belief that satisfy the properties of common knowledge of rationality and backward induction in the centipede game. Therefore, any other belief  $B_i$  would violate these properties. In this section, we do not focus on where this alternative belief is derived from. (It could be derived from a player's past game-play experience, for example. Since the population in general are not always rational. Even if someone is rational, he/she does make mistakes. All these factors can contribute to the forming of a player's belief about other players' behavior.) Instead, we want to characterize the equilibrium of the expanded game given this belief.

In summary, the game we are examining is as follows. Both players simultaneously select their beliefs before the start of the game. Once the belief is selected, it remains the same throughout the game. Given these beliefs regarding his opponent's behavior, they play the above centipede game. Each player's goal is to maximize his expected payoff given his chosen belief.

We explicitly impose that the beliefs will not be updated during the game. It simplifies the analysis so that we can emphasize the point we are trying to make. (Even if we allow for belief updating, we will not get back the SPE beliefs as long as the initial belief is somewhat incorrect.)

To analyze the modified centipede game, first note the following. If  $B_1$  is such that playing  $T$  at node 1 is the optimal action for player 1, then the game is over at node 1 no matter what belief player 1 has selected. The more interesting case is when  $T$  at node 1 is not the optimal action.

If player 1 chooses belief  $SPE_1$  and thus plays  $T$  at the first node, the game ends at the first node, with payoffs  $(a_1, b_1)$ . If player 1 chooses belief  $B_1$ , player 1 maximizes his expected payoff by choosing the node he plans to play  $T$ :

$$\begin{aligned} \max_{i \in \{1, 3, \dots, n-1\}} & p_2 a_2 + (1 - p_2) p_4 a_5 + \dots + (1 - p_2)(1 - p_4) \cdots (1 - p_{i-3}) p_{i-1} a_{i-1} \\ & + (1 - p_2)(1 - p_4) \cdots (1 - p_{i-3})(1 - p_{i-1}) a_i \end{aligned} \quad (1)$$

Let  $i = n_1^*$  denote an  $i$  that maximizes the above. (Note that there could be many such  $i$ 's that maximize the above.)

Consider player 2 at node 2. The optimal action with the belief of  $SPE_2$  is to end the game right away. In this case, the payoffs are  $(a_2, b_2)$ . If belief  $B_2$  is chosen, player 2 maximizes his expected payoff by choosing the node he plans to play  $T$ :

$$\begin{aligned} \max_{j \in \{2, 4, \dots, n\}} \quad & p_1 b_1 + (1 - p_1) p_3 b_3 + \dots + (1 - p_1)(1 - p_3) \cdots (1 - p_{j-3}) p_{j-1} b_{j-1} \\ & + (1 - p_1)(1 - p_3) \cdots (1 - p_{j-3})(1 - p_{j-1}) b_j \end{aligned} \quad (2)$$

Let  $j = n_2^*$  denote a  $j$  that maximizes the above. (Again, there could be many such  $j$ 's that maximize the above.)

The proposed pure strategy for player 1 is to select  $B_1$  and plan to play  $T$  at node  $n_1^*$ . The proposed pure strategy for player 2 is to select  $B_2$  and play  $P$  at node 2 (if player 1 played  $P$  at node 1), and plan to play  $T$  at node  $n_2^*$ . The game ends at node  $\min\{n_1^*, n_2^*\} \equiv n^*$ .

$\{B_1, n_1^*\}$  and  $\{B_2, n_2^*\}$  form a pure strategy “behavioral equilibrium” if player 1’s payoff is higher by selecting  $\{B_1, n_1^*\}$  than selecting  $\{SPE_1, 1\}$  given player 2’s strategy of playing  $T$  at node  $n_2^*$ , and player 2’s payoff is higher by selecting  $\{B_2, n_2^*\}$  than selecting  $\{SPE_2, 1\}$  given player 1’s strategy of playing  $T$  at node  $n_1^*$ . That is,

$$a_{n^*} \geq a_1, \text{ and } b_{n^*} \geq b_2.$$

In this behavioral equilibrium, players are better off selecting these non-SPE beliefs than selecting the SPE beliefs. Thus these beliefs are reinforced when the players play these games again later.

Now consider mixed strategy “behavioral equilibria”. Suppose that there are more than one  $j$ 's that maximize (2), or there are more than one  $i$ 's that maximize (1), mixed strategies could be used by the players. Let  $s_1 = (\dots, q_{i_1^*}, \dots, q_{i_2^*}, \dots, q_{i_k^*}, \dots)$  denote any of player 1’s optimal mixed strategies, where  $i_1^*, i_2^*, \dots, i_k^*$  are all of the numbers that maximizes (1). Similarly, let  $s_2 = (\dots, q_{j_1^*}, \dots, q_{j_2^*}, \dots, q_{j_k^*}, \dots)$  denote any of player 2’s optimal mixed strategies, where  $j_1^*, j_2^*, \dots, j_k^*$  are all of the numbers that maximizes (2). Then the outcomes of the game are determined by  $s_1$  and  $s_2$ .

$\{B_1, s_1^*\}$  and  $\{B_2, s_2^*\}$  form a mixed-strategy “behavioral equilibrium” if player 1’s payoff is higher by selecting  $\{B_1, s_1^*\}$  (comparing to  $\{SPE_1, 1\}$ ) given player 2’s strategy  $s_2^*$ , and player 2’s payoff is higher by selecting  $\{B_2, s_2^*\}$  (comparing to  $\{SPE_2, 2\}$ ) given player 1’s strategy  $s_1^*$ . Again, in this behavioral equilibrium, players are better off selecting these non-SPE beliefs than selecting those SPE beliefs.

Example 1 Consider the eight-move centipede game in Figure 2.

Suppose that  $B_1 = (0, 0, 0, 1)$  and  $B_2 = (0, 0, 0, 1)$ . Then it is straight-forward to obtain  $n_1^* = 7$ , and  $n_2^* = 6$ . That is, player 1 playing  $T$  at node 7 is optimal given  $B_1$ ,

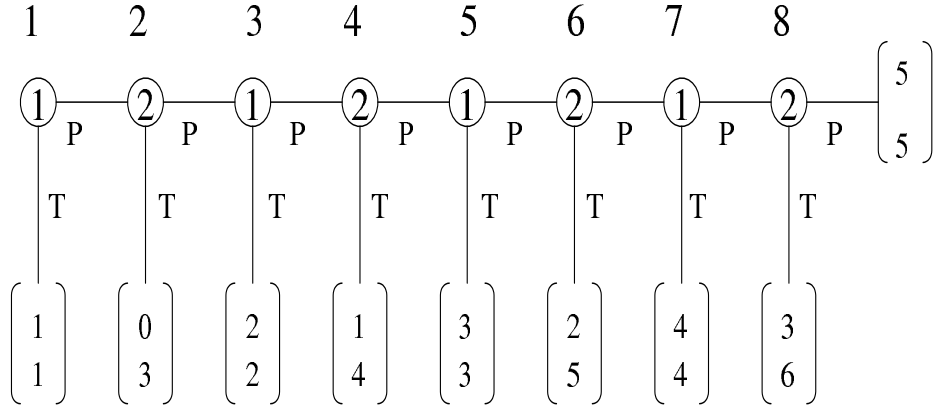


Figure 2: An Eight-Move Centipede Game

while player 2 playing  $T$  at node 6 is optimal given  $B_2$ . The minimum of  $n_1^*$  and  $n_2^*$ ,  $n^*$ , is 6; that is, the game ends at node 6, with payoffs (2,5).

It is easy to see that  $\{B_1, n_1^*\}$  and  $\{B_2, n_2^*\}$  form a behavioral equilibrium because  $a_{n^*} > a_1$ , and  $b_{n^*} > b_2$ .

Example 2 Consider the six-move centipede game in Figure 3.

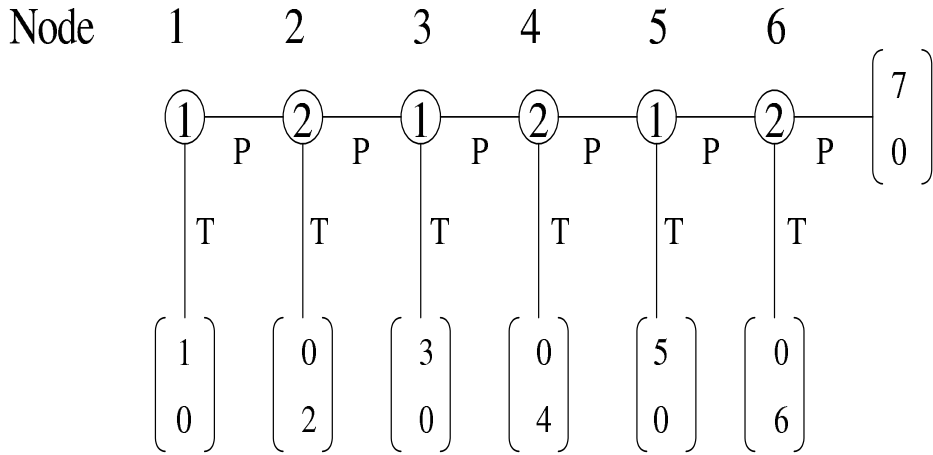


Figure 3: A Six-Move Centipede Game

In this game, we can construct pure-strategy behavioral equilibria similarly to the last example. Let  $B_1 = (0, 1, 0)$ , and  $B_2 = (0, 0, 1)$ . Then we have  $n_1^* = 3$ , and  $n_2^* = 4$ .



Therefore,  $n^* = \min\{n_1^*, n_2^*\} = 3$ ; that is, the game ends at node 3. This constitutes a behavioral equilibrium as the final outcome is (3,0), which is weakly better for both players than the SPE outcome of (1,0).

Now consider a mixed-strategy behavioral equilibrium. Suppose that  $B_1 = (0, p_4, 1)$  and  $B_2 = (0, p_3, 1)$ , with  $p_4 \in (0, 1)$  and  $p_3 \in (0, 1)$ . Given these beliefs, denote player 1's expected payoff of planning to play  $T$  at node  $i$  by  $E\Pi_1(i)$ . We have  $E\Pi_1(1) = 1$ ,  $E\Pi_1(3) = 3$ , and  $E\Pi_1(5) = p_4 \cdot 0 + (1 - p_4) \cdot 5$ . For player 1 to randomize between playing  $T$  at node 3 and playing  $T$  at node 5, we should set  $E\Pi_1(3) = E\Pi_1(5)$ ; that is,  $p_4 = \frac{2}{5}$ .

Similarly, for player 2,  $E\Pi_2(2) = 2$ ,  $E\Pi_2(4) = p_3 \cdot 0 + (1 - p_3) \cdot 4$ , and  $E\Pi_2(6) = 0$ . Suppose that  $p_3 < \frac{1}{2}$ . Then  $n_2^* = 4$ .

To construct a behavioral equilibrium, player 1's mixed strategy  $(0, q_3, 1)$  must satisfy the following two conditions regarding the each player's actual payoffs. First, for player 1,  $q_3 \cdot 3 + (1 - q_3) \cdot 0$  is at least 1, which is player 1's payoff by following SPE strategy and playing  $T$  at node 1. This gives us  $q_3 \geq \frac{1}{3}$ . Second, for player 2,  $q_3 \cdot 0 + (1 - q_3) \cdot 4$  must be at least 2, which is player 2's payoff by following SPE strategy and playing  $T$  at node 2. This gives us  $q_3 \leq \frac{1}{2}$ . Therefore, any  $q_3 \in [\frac{1}{3}, \frac{1}{2}]$  would satisfy these two conditions.

To summarize,  $B_1 = (0, \frac{2}{5}, 1)$ ,  $s_1 = (0, q_3, 1)$ ,  $B_2 = (0, p_3, 1)$ ,  $s_2 = (0, 1, 1)$ , where  $q_3 \in [\frac{1}{3}, \frac{1}{2}]$ , and  $p_3 \in [0, \frac{1}{2}]$  form a mixed-strategy behavioral equilibrium.

### 3 Analyzing Previous Centipede Game Experiments

McKelvey and Palfrey (1992) report the results of seven different sessions of the centipede game experiments. Sessions 1 to 3 are four-move centipede games with the following payoffs:  $(a_1, b_1) = (0.4, 0.1)$ ,  $(a_2, b_2) = (0.2, 0.8)$ ,  $(a_3, b_3) = (1.6, 0.4)$ ,  $(a_4, b_4) = (0.8, 3.2)$ , and  $(a_5, b_5) = (6.4, 1.6)$ . Session 4 is a high-payoff four-move centipede game where the payoffs are quadruppled. Sessions 5 to 7 are six-move centipede games with the following payoffs:  $(a_1, b_1) = (0.4, 0.1)$ ,  $(a_2, b_2) = (0.2, 0.8)$ ,  $(a_3, b_3) = (1.6, 0.4)$ ,  $(a_4, b_4) = (0.8, 3.2)$ ,  $(a_5, b_5) = (6.4, 1.6)$ ,  $(a_6, b_6) = (3.2, 12.8)$ , and  $(a_7, b_7) = (25.6, 6.4)$ .

Table IIA in McKelvey and Palfrey (1992) describes the proportion of observations at each terminal node. In that table,  $f_i$  is used to denote the proportion of games that ends at node  $i$ . From these  $f_i$ 's, we can calculate a player's strategy as follows. For the four-move game, let  $q_1$  and  $q_3$  be the proportion of player 1 who plan to choose TAKE at node 1 and at node 3 respectively. (Therefore, the proportion of player 1 choosing Pass at node 3 is equal to  $1 - q_1 - q_3$ .) Similarly, let  $q_2$  and  $q_4$  be the proportion of player 2 who plan to choose TAKE at node 2 and at node 4 respectively, and thus the

proportion of player 2 choosing Pass at node 4 is equal to  $1 - q_2 - q_4$ . Then  $q_1 = f_1$ ,  $(1 - q_1)q_2 = f_2$ ,  $(1 - q_2)q_3 = f_3$ , and  $(1 - q_1 - q_3)q_4 = f_4$ . We define  $q_i$  similarly in the six-move game. Then we have  $q_1 = f_1$ ,  $(1 - q_1)q_2 = f_2$ ,  $(1 - q_2)q_3 = f_3$ ,  $(1 - q_1 - q_3)q_4 = f_4$ ,  $(1 - q_2 - q_4)q_5 = f_5$ , and  $(1 - q_1 - q_3 - q_5)q_6 = f_6$ . The results are reported in the following table.

Table 1: Players' Strategies and Optimal Actions

Session		$q_1$	$q_2$	$q_3$	$q_4$	$q_5$	$q_6$	Optimal Action
1	player 1	.06		.61				Take at Node 3 (61%)
	player 2		.28		.61			Take at Node 4 (61%)
2	player 1	.10		.69				Take at Node 3 (69%)
	player 2		.42		.52			Take at Node 4 (52%)
3	player 1	.06		.52				Pass at Node 3 (42%)
	player 2		.46		.33			Take at Node 4 (33%)
4	player 1	.15		.57				Take at Node 3 (57%)
	player 2		.44		.39			Take at Node 2 (44%)
5	player 1	.02		.43		.50		Take at Node 5 (50%)
	player 2		.09		.51		.20	Take at Node 4 (51%)
6	player 1	.00		.04		.70		Take at Node 5 (70%)
	player 2		.02		.48		.42	Take at Node 4 (48%)
7	player 1	.00		.15		.55		Take at Node 5 (55%)
	player 2		.07		.51		.40	Take at Node 4 (51%)

It is hard to infer a player's belief in playing these games, since many different beliefs could lead to the same observed strategy. Therefore, in each session, we take a player's rivals' revealed strategies as the player's belief and calculate the player's optimal action according to that belief. In the calculations, we assign the players a utility function with constant degree of absolute risk aversion of 0.5. That is,  $U_i(x) = -e^{-0.5x}$ , where  $x$  is the amount of money earned in one game. Therefore, the players are modestly risk averse. The results are reported in Table 1 as well. The percentage number after each optimal action is the percentage of players actually choosing the implied optimal action in that session. As we can see from the table, the majority of the players chose the implied optimal action in all but session 3. It suggests that the behavior of the majority of the players can be explained by our theory.

## 4 Concluding Remarks

In this paper, we propose a concept of behavioral equilibrium explaining the behavior of players in centipede games. Players' behavior is usually different from what game theory

has predicted. We allow players to abandon the logic of subgame perfect equilibrium and choose a belief that is formed (e.g. from their previous experience in the situation). Under certain conditions, the players are better off by abandoning the subgame perfect equilibrium belief and choose the alternative belief instead. This reinforces the players' subjective opinion that subgame perfect equilibrium may not work well in these games. Hence, alternative beliefs become the beliefs of choice. We support our theory by re-examining some previous centipede game experiments.

## References

- [1] Aumann, Robert, "Irrationality in Game Theory", in Dasgupta, Partha, et al, eds. *Economic Analysis of Markets and Games: Essays in Honor of Frank Hahn*, MIT Press: Cambridge and London, 1992, pp.214-27.
- [2] Basu, Kaushik, "Strategic Irrationality in Extensive Games", *Mathematical Social Sciences*, Vol. 15 (1988), pp. 247-60.
- [3] Caplan, Bryan, "Rational Irrationality and the Microfoundations of Political Failure", *Public Choice*, Vol. 107 (2001), pp. 311-31.
- [4] Epstein, Seymour, et al, "Irrational Reactions to Negative Outcomes: Evidence for Two Conceptual Systems", *Journal of Personality and Social Psychology*, Vol. 62, 328-39, 1992.
- [5] Goeree, Jacob, and Charles Holt, "Ten Little Treasures of Game Theory and Ten Intuitive Contradictions", *American Economic Review*, Vol. 91 (2001), pp. 1402-22.
- [6] McKelvey, Richard, and Thomas Palfrey, "An Experimental Study of the Centipede Game", *Econometrica*, Vol. 60, pp. 803-36, 1992.
- [7] McKelvey, Richard, and Thomas Palfrey, "Quantal Response Equilibria for Extensive Form Games", *Experimental Economics*, Vol.1, pp.9-41, 1998.
- [8] Rosenthal, Robert, "Games of Perfect Information, Predatory Pricing, and the Chain Store Paradox", *Journal of Economic Theory*, Vol.25, pp.92-100, 1982.
- [9] Tirole, Jean, "Rational Irrationality: Some Economics of Self- Management", *European Economic Review*, Vol. 46 (2002), pp. 633-55.
- [10] Zauner, Klaus, "A Payoff Uncertainty Explanation of Results in Experimental Centipede Games," *Games and Economic Behavior*, Vol. 26 (1999), pp. 157-185.