

# Evolutionary Stability of Prospect Theory Preferences

Marc Oliver Rieger\*

February 25, 2009

## Abstract

We demonstrate that in simple  $2 \times 2$  games (cumulative) prospect theory preferences can be evolutionarily stable, i.e. a population of players with prospect theory preferences can not be invaded by more rational players. This holds also if probability weighting is applied to the probabilities of mixed strategies. We also show that in a typical game with infinitely many strategies, the “war of attrition”, probability weighting is evolutionarily stable. Our results may help to explain why probability weighting is generally observed in humans, although it is not optimal in usual decision problems.

**Keywords:** prospect theory, existence of Nash equilibria, evolutionary stability.

**JEL classification:** C70, C73, D81.

## 1 Introduction

We study the influence of prospect theory preferences on the outcome of two player games. We focus on the effect of probability weighting on the probabilities for mixed strategies. A priori one might assume that probability weighting reduces the (rational) payoffs<sup>1</sup> a player receives in a game, since it leads to irrational decisions. In this article, however, we demonstrate that there are many situations where probability weighting of the players is evolutionarily stable, in particular in a class of simple  $2 \times 2$  games related to matching pennies games which we call *social control games* (see Sec. 2.1-2.2) and in the “war of attrition” (Sec. 2.3).

We suggest that our results provide a possible explanation for the “probability weighting puzzle”, i.e. the question why humans tend to overweight small probabilities, given that this leads to suboptimal decisions (as compared to the expected utility benchmark): when considering interactions between individuals, the seemingly irrational probability weighting can become advantageous and evolutionarily

---

\*University of Zürich, ISB, Plattenstrasse 32, 8032 Zürich, Switzerland, rieger@isb.uzh.ch and University of Bielefeld, IMW.

<sup>1</sup>based on linear probability weighting according to subjective expected utility theory

stable. Since humans do not only face simple (single person) decision problems, but manifold interactions with others, on average a neutral probability weighting is usually not optimal (Sec. 3).

### 1.1 Expected utility theory and prospect theory

Since its origin (v. Neumann 1928, von Neumann & Morgenstern 1944) game theory has been closely connected to expected utility theory. The role of this decision model as a normative theory is uncontroversial. Recent years, however, have seen substantial progress on the understanding of differences between the actual, sometimes irrational decisions of individuals, and rational decisions according to the expected utility theory. There are now several models to describe decisions under risk, in particular prospect theory, as developed by Kahneman & Tversky (1979) and Tversky & Kahneman (1992), a model for which Daniel Kahneman was awarded with the Nobel Prize for economics in 2002. Since prospect theory and its variants are nowadays the most frequently used behavioral decision models, we concentrate on them.<sup>2</sup>

Prospect theory modifies classical expected utility theory in several ways:

1. Unlike expected utility theory, not the final wealth is evaluated, but the pay-offs are framed as gains and losses with respect to a reference point; they are called “prospects”.
2. Losses loom larger than gains, hence the marginal utility in losses is larger than in gains.
3. Small probabilities are overweighted and moderate to large probabilities are underweighted.

Mathematically, the first two features are reflected in a two-part S-shaped *value function* (which replaces the usual utility function) – concave in gains and convex in losses. The prototypical example has been given in Tversky & Kahneman (1992) for  $\alpha, \beta \in (0, 1)$  and  $\lambda > 1$ :

$$u(x) := \begin{cases} x^\alpha, & x \geq 0 \\ -\lambda(-x)^\beta, & x < 0. \end{cases} \quad (1)$$

The third feature is captured by weighting the probability distribution by an S-shaped function, the so-called *weighting function*  $w$ . The original example of Tversky & Kahneman (1992) is given by

$$w(F) := \frac{F^\gamma}{(F^\gamma + (1 - F)^\gamma)^{1/\gamma}}. \quad (2)$$

---

<sup>2</sup>We mention approaches using the Choquet integral that are related to cumulative prospect theory, see (Gilboa & Schmeidler 1992) and (Schmeidler 1989), that could be treated in a similar matter.

In the classical form of prospect theory (PT), the function  $w$  is applied directly to the probabilities of the different outcomes, resulting in an overweighting of small probabilities, regardless of their associated outcome. The generalization of this form for outcomes  $x_i$  with probabilities  $p_i$  is given by

$$PT(p) = \sum_{i=1}^n w(p_i)u(x_i), \quad (3)$$

where  $PT(p)$  is the *subjective utility* of a probability distribution  $p$  (Kahneman & Tversky 1979, Schneider & Lopes 1986, Wakker 1989).

The updated version of prospect theory, called *cumulative prospect theory* (short: CPT), weights cumulative probabilities  $F_i := \sum_{j=1}^i p_j$ , where outcomes are ordered by their payoffs, and the weight factor for the  $i$ -th outcome is  $w(F_i) - w(F_{i-1})$ .<sup>3</sup> The result is that only low probability events *with extreme outcomes* are overweighted. CPT helped in recent years to explain various effects in decision theory, economics and finance.

## 1.2 Prospect theory preferences in games

One of the most interesting effects of prospect theory on the analysis of games is the interplay between probability weighting and mixed strategies. This effect makes it necessary to extend the previous work on non-expected utility preferences in games. Let us consider a finite normal-form game with two players (without chance moves). In this game, a player  $i$  can choose from the strategy set  $S_i$ ,  $i \in \{A, B\}$  of (finitely many) pure strategies. We denote the set of all combinations of pure strategies  $S := \times_{i \in \{A, B\}} S_i$ . The set of probability measures on  $S_i$  is denoted by  $M_i$  and describes the mixed strategies of player  $i$ . The combinations of mixed strategies are denoted by  $M := \times_{i \in \{A, B\}} M_i$ . The payoff (in utility units) of the game for the  $i$ -th player is given by  $u_i: S \rightarrow \mathbb{R}$ . The game can then be written as  $(S_i, u_i)_{i \in \{A, B\}}$ .

The total utility  $U$  that a player, say player A, obtains for some mixed strategy play  $m = (m_1, \dots, m_n) \in M$  depends on the underlying decision model. In the case of EUT, this utility becomes

$$U_A^{EUT}(m) = \sum_{s=(s_A, s_B) \in S} m_A(s_A)m_B(s_B)u_A(s),$$

where  $m_i(s)$  is the probability of player  $i$  to play strategy  $s$ .

We consider now probability weighting functions  $w_i$ . If the player weights the probabilities with which the other player chooses his mixed strategies using this probability weighting function, we obtain for the utility of the first player

$$U_A^{PT}(m) = \sum_{s=(s_A, s_B) \in S} m_A(s_A)w_A(m_B(s_B))u_A(s). \quad (4)$$

<sup>3</sup>For a precise formula see Tversky & Kahneman (1992).

We assume that the players do not weight the probabilities of their own strategies<sup>4</sup>. Here and in the remaining part of this article we assume that the reference point of the value function  $u$  is fixed, see Rieger & Koch (2009) for generalizations. In the case of cumulative prospect theory (CPT), we need to rank the possible outcomes, before we can compute the probability weighting. To simplify notation we denote the potential outcomes for player A by  $u_A(i, k)$ , where  $i$  is his own strategy and  $k$  is the strategy played by player B. In order to define cumulative probabilities, we sort these outcomes first, therefore let us define permutations  $\sigma_i^A$  on  $\{1, \dots, n\}$  such that

$$u_A(i, \sigma_i^A(k)) \leq u_A(i, \sigma_i^A(k+1)), \quad \text{for all } k = 1, \dots, n-1.$$

Now we can define the cumulative probabilities<sup>5</sup> of player B's actions as seen from player A by

$$\begin{aligned} F_A(i, k) &:= \sum_{l=1}^k m_{\sigma_i^A(l)}^B, \\ F_A(i, 0) &:= 0, \end{aligned}$$

and the CPT-utility becomes

$$U_A^{CPT} = \sum_{i=1}^n \left( m_i^A \sum_{j=1}^n (w_A(F_A(i, j)) - w_A(F_A(i, j-1))) u_A(i, \sigma_i^A(j)) \right).$$

This is essentially the form used by Goeree, Holt & Pfafrey (2003) in the special case of  $2 \times 2$  games. We will later show how to generalize this formulation to a game with infinitely many pure strategies.

Nash equilibria in this setting can be defined as usual:

**Definition 1.1.** *We call a strategy  $\hat{m} \in M$  a mixed PT-Nash equilibrium if for all  $i = 1, \dots, n$  and all  $m \in M$  with  $m_k = \hat{m}_k$  for  $k \neq i$  we have  $U_i(\hat{m}) \geq U_i(m)$ , where  $U_i = U_i^{PT}$  is given by (4).*

*Analogously, we say that  $\hat{m} \in M$  is a mixed CPT-Nash equilibrium if for  $i = 1, 2$  and all  $m \in M$  with  $m_k = \hat{m}_k$  for  $k \neq i$  we have  $U_i(\hat{m}) \geq U_i(m)$ , where  $U_i = U_i^{CPT}$  is given by (5).*

The existence of Nash equilibria in the finite game case has been proven in Rieger & Koch (2009) under the assumption of fixed (or continuous) reference points. Extensions of the equilibria concept are also introduced

<sup>4</sup>There is an older approach by Dekel, Safra & Segal (1991) for non-expected utility theory which weights also the probabilities of the player's own strategies. This approach cannot be extended to cumulative prospect theory, since it is not possible to rank both the player's and the opponent's strategies simultaneously by the payoff. There are also conceptual reasons in favor of the approach used here, see Rieger & Koch (2009).

<sup>5</sup>There are slight differences in the precise definition of CPT in the literature. In the original formation (Tversky & Kahneman 1992), cumulative probabilities have been used in losses, but decumulative probabilities in gains. For our analysis, this difference would only be quantitative, but does not change the qualitative results: one could in any case always adjust the probability weighting function accordingly.

## 2 Evolutionary stability of probability weighting

Given that prospect theory is not a normative theory, but rather describes systematic deviations from rational choices, it has to be expected that players who show such deviations will fare worse in games than players with rational preferences. We will see, however, that there are simple classes of games where probability weighting of the players leads in fact to an increase in the players' payoffs and is even evolutionarily stable! We motivate this first with a simple example (Sec. 2.1) before we generalize this result (Sec. 2.2). Finally, we show that such results can also be obtained for games with infinitely many strategies (Sec. 2.3).

### 2.1 A simple $2 \times 2$ game

We consider the following game that has a similar structure as the classical matching pennies game:

$$\begin{array}{c|cc}
 & \text{Player B} & \\
 \hline
 \text{Player A} & (4, 2) & (-2, 3) \\
 & (3, 2) & (0, 0)
 \end{array} \tag{5}$$

We denote a mixed strategy of player A by  $p \in [0, 1]$ , meaning the probability with which he chooses his first strategy. Accordingly, we denote a mixed strategy of player B by  $q \in [0, 1]$ . Without probability weighting this game has the unique mixed Nash equilibrium  $(p, q) = (2/3, 2/3)$ , as a short computation shows.

We now use the standard form of probability weighting (2) and denote the probability weighting parameter of player A and B by  $\alpha$  and  $\beta$ , respectively.

The PT-utilities for the players are then

$$\begin{aligned}
 U_A(p, q) &= 4pw_\alpha(q) - 2pw_\alpha(1 - q) + 3(1 - p)w_\alpha(q), \\
 U_B(p, q) &= 2w_\beta(p)q + 2w_\beta(1 - p)q + 3w_\beta(p)(1 - q).
 \end{aligned}$$

To compute the CPT-utilities we first sort the outcomes for each player given that he plays a certain strategy by their payoffs and obtain

$$\begin{aligned}
 U_A(p, q) &= 4p(1 - w_\alpha(1 - q)) - 2pw_\alpha(1 - q) + 3(1 - p)(1 - w_\alpha(1 - q)), \\
 U_B(p, q) &= 2q + 3(1 - w_\beta(1 - p))(1 - q).
 \end{aligned}$$

Probability weighting is a behavioral bias and hence leads usually to suboptimal decisions. The key observation of this article, however, is that in certain games this behavioral bias can be evolutionarily stable, and hence it can be optimal to have prospect theory preferences, rather than expected utility preferences. To make this question more precise: We follow the ideas of Güth & Yaari (1992) and Ely & Yilankaya (2001) and consider a “meta-game” where we assume that in the single game always the Nash equilibrium is played (with changing roles of the players) and that on the meta-level the probability weighting can change. This

could be interpreted as two time-scales of adaptation: the direct response to a game is adapting quickly with respect to the opponents' responses, whereas the overall behavioral patterns are either inherited or part of a cultural setting and therefore can change only very slowly.

We also introduce a “one-sided” form of evolutionary stability motivated by the fact that we are mostly interested to know whether a strong overweighting is stable against invasion of weaker overweighting or rational behavior (i.e. no overweighting). The reason behind this is that our model ignores the fact that probability weighting does not only have an impact on the specific game under study, but also on other situations an individual may face. In particular, there are instances when an individual has to make a decision without interaction with other players. Here, it seems clear that a rational decision procedure, i.e. no probability weighting, would be evolutionarily optimal. Both effects together would then lead to an evolutionarily stable amount of probability weighting. For simplicity, and since this is the interesting part of the result, we focus on the question whether a certain degree of probability weighting can be evolutionarily stable *against a lower degree of probability weighting* when considering social control games. We call this variant of evolutionary stability “semi-stability”.

We give the following definitions:

**Definition 2.1.** We call an individual with a probability weighting  $\gamma \in (0, 1]$  a  $\gamma$ -weighter.

We denote the rational utility that a  $\gamma$ -weighter obtains when playing the game as player A against a  $\delta$ -weighter by  $U_A(\gamma, \delta)$  and the rational utility that a  $\gamma$ -weighter obtains when playing the game as player B against a  $\delta$ -weighter by  $U_B(\delta, \gamma)$ . Define the average utility by  $U(\gamma, \delta) := (U_A(\gamma, \delta) + U_B(\delta, \gamma))/2$ .

A probability weighting  $\gamma \in (0, 1)$  is called evolutionarily semi-stable if for all  $\delta \in (0, 1]$  with  $\delta > \gamma$  and for all sufficiently small  $\varepsilon > 0$  the rational utility of  $\gamma$ -weighters is larger than the rational utility of  $\delta$ -weighters, where the proportions of  $\gamma$ - and  $\delta$ -weighters are  $1 - \varepsilon$  and  $\varepsilon$ , respectively, i.e.

$$\varepsilon U(\gamma, \delta) + (1 - \varepsilon)U(\gamma, \gamma) > \varepsilon U(\delta, \delta) + (1 - \varepsilon)U(\delta, \gamma). \quad (6)$$

A probability weighting  $\gamma \in (0, 1)$  is called evolutionarily stable if this holds even for  $\delta < \gamma$ .

We can now state the evolutionary (semi-) stability of prospect theory preferences in the case of game (5):

**Proposition 2.2.** In the game (5) probability weighting is evolutionarily semi-stable in both the PT and the CPT setting, i.e. for sufficiently large  $\gamma < 1$  a population of  $\gamma$ -weighters can not be invaded by  $\delta$ -weighters with  $\delta > \gamma$ , i.e. individuals with a smaller (or none) behavioral bias.

In the CPT setting there even exists a  $\bar{\gamma} \approx 0.75$  that is evolutionarily stable, i.e. a population of  $\gamma$ -weighters can not be invaded by  $\delta$ -weighters with  $\delta \neq \gamma$ .

This result implies that the seemingly irrational prospect theory preferences can be “rational” in games – rational on the meta-level of the evolution of preferences, in the sense of evolutionary stability.

*Proof.* The results for PT will be a special case of the general evolutionary semi-stability result in the next section.

To prove evolutionary stability in the CPT case we first compute the (unique) Nash equilibrium in this case: The first order condition is

$$0 = \frac{d}{dp} U_A(p, q) = 1 - 3w_\alpha(1 - q).$$

Inverting the weighting function we obtain  $q = 1 - w_\alpha^{-1}(1/3)$ . Similarly, we obtain  $p = 1 - w_\beta^{-1}(1/3)$ .

For the standard form of the weighting function (2) there exists a  $\bar{\gamma} \in (0, 1)$  that minimizes  $w_\gamma^{-1}(1/3)$ , as can be seen from Fig. 1.<sup>6</sup> We claim that  $\bar{\gamma}$  (which is approximately 0.75) is evolutionarily stable.

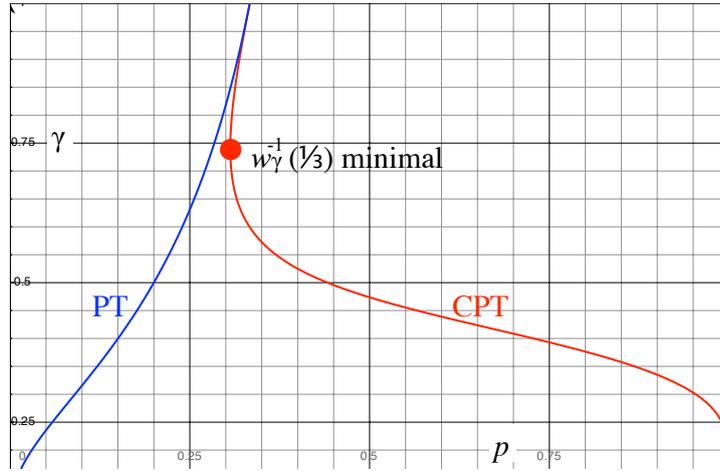


Figure 1: The minimal value for  $w_\gamma^{-1}(1/3)$  is approximately  $\bar{\gamma} = 0.75$  for CPT and zero for PT, in both cases different from the rational behavior  $\gamma = 1$ . Hence, in the CPT model the evolutionarily stable strategy is to weight probabilities with a probability weighting parameter of  $\bar{\gamma} \approx 0.75$ . This value changes slightly when changing the parameters of the game.

Define  $p_{\gamma\delta}$  as player A’s strategy in the Nash equilibrium when player A has a probability weighting  $\gamma$  and player B has a probability weighting  $\delta$ . As computed

<sup>6</sup>The weighting function in (2) can be inverted for all  $\gamma > 0.3$  which is sufficient in our case.

above, we have particularly  $p_{\alpha\beta} = q_{\alpha\beta}$  for all  $\alpha, \beta \in (0, 1]$ . We compute

$$\begin{aligned} U_A(\gamma, \delta) &:= \mathbb{E}U_A(p_{\gamma\delta}, q_{\gamma\delta}) \\ &= 3p_{\gamma\delta}q_{\gamma\delta} + 3q_{\gamma\delta} - 2p_{\gamma\delta} \\ U_B(\delta, \gamma) &:= \mathbb{E}U_B(p_{\delta\gamma}, q_{\delta\gamma}) \\ &= -3p_{\delta\gamma}q_{\delta\gamma} + 3p_{\delta\gamma} + 2q_{\delta\gamma} \end{aligned}$$

and

$$\begin{aligned} U(\gamma, \delta) &= \frac{1}{2}(U_A(\gamma, \delta) + U_B(\delta, \gamma)) \\ &= 3q_{\gamma\delta}. \end{aligned}$$

We can now show evolutionary stability using that we had chosen  $\bar{\gamma}$  such that  $w_{\bar{\gamma}}^{-1}(1/3) < w_{\delta}^{-1}(1/3)$  for all  $\delta \neq \bar{\gamma}$ :

$$\begin{aligned} \varepsilon U(\bar{\gamma}, \delta) + (1 - \varepsilon)U(\bar{\gamma}, \bar{\gamma}) &= 3\varepsilon q_{\bar{\gamma}\delta} + 3(1 - \varepsilon)p_{\bar{\gamma}\bar{\gamma}} \\ &= 3\varepsilon(1 - w_{\bar{\gamma}}^{-1}(1/3)) + 3(1 - \varepsilon)(1 - w_{\bar{\gamma}}^{-1}(1/3)) \\ &> 3\varepsilon(1 - w_{\delta}^{-1}(1/3)) + 3(1 - \varepsilon)(1 - w_{\delta}^{-1}(1/3)) \\ &= 3\varepsilon p_{\delta\delta} + 3(1 - \varepsilon)p_{\bar{\gamma}\delta} \\ &= \varepsilon U(\delta, \delta) + (1 - \varepsilon)U(\delta, \bar{\gamma}). \end{aligned}$$

Thus  $\bar{\gamma} \approx 0.75$  is the only evolutionarily stable amount of probability weighting.  $\square$

## 2.2 A general class of $2 \times 2$ games where probability weighting is evolutionarily stable

Is the game presented in the previous section a somehow “pathological” case or are there meaningful classes of games where probability weighting is evolutionarily stable? In this section we will study a relatively broad class of games which we call “social control games” that share the same features with the preceding example (5). We will give a motivation for this class of games below. The class shares certain features with the matching pennies game, but has less symmetry, in particular it does not contain zero-sum games.

Afterwards we will study another (unrelated) class of games, the “war of attrition”, that also show evolutionary stability of probability weighting. Thus the phenomenon seems to be widespread.

In both cases we will present for simplicity only the prospect theory case.

Let us consider general  $2 \times 2$  games given by the following payoff matrix, given in utility values:

	$q$	$1 - q$
$p$	$(A_1, B_1)$	$(A_3, B_3)$
$1 - p$	$(A_2, B_2)$	$(A_4, B_4)$

To ease computations we normalize the individual utilities such that  $A_4 = B_4 = 0$  by adding a fixed amount to *all* payoffs of a player.<sup>7</sup>

Let us assume, B plays the mixed strategy  $q$ , then the subjective utility for player A, given that he plays his first strategy (in other words  $p = 1$ ) will be

$$U_A(1, q) = w_\alpha(q)A_1 + w_\alpha(1 - q)A_3.$$

Analogously, the utility when playing his second strategy will be

$$U_A(0, q) = w_\alpha(q)A_2.$$

(Remember that  $A_4 = 0$ .) Overall, his PT utility when he plays a mixed strategy  $p$  is

$$U_A(p, q) = pw_\alpha(q)A_1 + pw_\alpha(1 - q)A_3 + (1 - p)w_\alpha(q)A_2.$$

The utility for player B is accordingly

$$U_B(p, q) = w_\beta(p)qB_1 + w_\beta(1 - p)qB_2 + w_\beta(p)(1 - q)B_3.$$

It is now easy to compute the mixed Nash equilibria of such a game (if at least one non-pure Nash equilibrium exists). The result of this standard computation is as follows: Assume that  $A_1 \neq A_2$  and  $B_1 \neq B_3$  and that

$$A_0 := \frac{A_3}{A_2 - A_1} > 0, \quad B_0 := \frac{B_2}{B_3 - B_1} > 0.$$

Then there exists a (unique) mixed Nash equilibrium  $(p, q)$  with

$$p = \frac{B_0^{1/\beta}}{1 + B_0^{1/\beta}}, \quad q = \frac{A_0^{1/\alpha}}{1 + A_0^{1/\alpha}}.$$

We see that varying the values of  $\alpha$  and  $\beta$  will also shift the position of the mixed Nash equilibrium  $(p, q)$ . This shift is monotone and its direction depends on the numbers  $A_0$  and  $B_0$ . More precisely, a short computation gives the following result:

**Lemma 2.3.** *The dependence of  $p$  and  $q$  on  $\alpha$  and  $\beta$  is as follows:*

1. If  $A_0 \in (0, 1)$ , decreasing  $\alpha$  decreases  $q$ .
2. If  $A_0 > 1$ , decreasing  $\alpha$  increases  $q$ .
3. If  $B_0 \in (0, 1)$ , decreasing  $\alpha$  decreases  $p$ .
4. If  $B_0 > 1$ , decreasing  $\alpha$  increases  $p$ .

---

<sup>7</sup>Since we later discuss the total utility gained from certain iterative plays, we cannot normalize further by adding a fixed number to all entries, say, of a column for player A.

*In the limit cases of  $\alpha \rightarrow 0$  and  $\beta \rightarrow 0$ ,  $p$  and  $q$  converge to 0 or 1. In the limit of  $\alpha \rightarrow 1$  and  $\beta \rightarrow 1$ ,  $p$  and  $q$  converge to  $\bar{p}$  and  $\bar{q}$ , the Nash equation of the game without probability weighting.*

Such a shift of the mixed Nash equilibrium in the case of “matching pennies games” has been observed already in Goeree et al. (2003, page 15 ff.).

How do these results change if we replace PT by CPT? It is quite obvious from the analysis above that the changes will be minor: the formulae for the subjective utilities will slightly differ and therefore it will in general not be possible to find a nice closed form for the mixed CTP-Nash equilibria. The qualitative behavior, however, remains similar, and in particular the essential structure of Lemma 2.3 will carry over to this situation as well besides that the convergence for  $\alpha, \beta \rightarrow 0$  might differ, as the example of game (5) demonstrates. In the remaining part of this article, we will not consider this, and instead only focus on PT.

In the following we concentrate on a special case of  $2 \times 2$  games that describe a certain kind of *social interaction*. Behavior that is directed towards the common wealth and not towards selfish goals can either be mutually enforced on the level of iterated strategies (like in the famous iterated prisoner’s dilemma) or via direct control by others. Here we focus on the latter case, where we will show that probability weighting can improve this control mechanism. We will call such games “social control games”. In these games, one of the players can enforce some social norm and the other player can either follow the norm or deviate. Deviation from the social norm would increase the utility for the deviator and decrease the utility of the enforcer, but the enforcer can check the deviator’s behavior. If he notices a deviation, he punishes the deviator. On the other hand, the enforcer would not like to check too much, since controlling is costly for him. We also assume that there is no positive effect on the controller if he checks and catches a deviator and that a non-deviating player does not profit when the controller checks on him.

Variants of this game occur naturally in all societies. Examples in our cultural context could include the interaction between employees (who decide between working and being lazy) and their employer (who can check on them), or the interaction between students (who can study topics for an exam or skip them) and their professor (who can check on some of the topics in the exam).

Formally, we can define this class of games by imposing conditions on the general game.

	follow	deviate
do not check	$(A_1, B_1)$	$(A_3, B_3)$
check	$(A_2, B_2)$	$(0, 0)$

The basic condition is that a social norm should correspond to a higher common wealth (defined as the sum of the players’ utilities), therefore we assume that  $A_1 + B_1 > \max\{A_2 + B_2, A_3 + B_3, 0\}$ . From this condition it is clear that the strategic pair optimizing the “common wealth” is for A not to check and for B to follow.

Social control and the possibility of deviating from it can be expressed by additional conditions:

- checking should not be for free, i.e.  $A_1 > A_2$  (otherwise, checking would be done routinely);
- deviating when not checked should be desirable to B, e.g. the employee (otherwise no need for control!), i.e.  $B_3 > B_1$ ;
- catching a deviator enforces the norm and is therefore better for A, e.g. the employer, than not catching a deviator, i.e.  $0 > A_3$ ;
- when checked, it is better not to deviate, i.e.  $B_2 > 0$ ;
- catching a deviator is still worse than if he had followed the norm at the first place, i.e.  $A_2 > 0$ ;
- there is no “honesty premium”, i.e.  $B_1 \geq B_2$ .

Under these conditions, the only Nash equilibrium is a mixed strategy where A checks with a probability  $\bar{p}$  and B follows with probability  $\bar{q}$ . In the rational case, the probabilities of this Nash equilibrium are

$$\bar{p} = \frac{B_0}{1 + B_0}, \quad \bar{q} = \frac{A_0}{1 + A_0},$$

where

$$A_0 := \frac{A_3}{A_2 - A_1}, \quad B_0 := \frac{B_2}{B_3 - B_1}.$$

A social norm should be followed at least more than half of the time and should be accepted enough in order to be checked on in less than half of the time, otherwise it would more be an exception rather than a norm. Therefore we assume that  $\bar{p}, \bar{q} > 0.5$ , in other words, we assume that

$$A_0 = \frac{A_3}{A_2 - A_1}, \quad B_0 = \frac{B_2}{B_3 - B_1} > 1,$$

as can be seen by a small computation.

We call a game satisfying all of these properties a *social control game* as summarized in the following definition:

**Definition 2.4.** We call a game with the payoff matrix

$$\begin{array}{cc} (A_1, B_1) & (A_3, B_3) \\ (A_2, B_2) & (0, 0) \end{array}$$

that satisfies

$$(i) \quad A_1 + B_1 > \max\{A_2 + B_2, A_3 + B_3, 0\},$$

(ii)  $A_1 > A_2$ ,  $B_3 > B_1$ ,  $0 > A_3$  and  $B_2 > 0$ ,

(iii)  $\frac{A_3}{A_2 - A_1} > 1$ ,  $\frac{B_2}{B_3 - B_1} > 1$

(iv)  $A_2 > 0$ ,

(v)  $B_1 \geq B_2$

a social control game.

How does the Nash equilibrium of a social control game change if one or both of the players overweight small probabilities? If we look at the PT model (4), the optimal strategies  $p$  and  $q$  are given by

$$p = \frac{B_0^{1/\beta}}{1 + B_0^{1/\beta}}, \quad q = \frac{A_0^{1/\alpha}}{1 + A_0^{1/\alpha}}.$$

Let us now have a look on how the mixed strategies are shifted in social control games. By assumption (iii),  $A_0 > 1$  and  $B_0 > 1$ , therefore by Lemma 2.3 both  $p$  and  $q$  increase when  $\alpha$  and  $\beta$  decrease, i.e. in the case of stronger overweighting. (Remember that assumption (ii) ensures that the Nash equilibrium is mixed, otherwise the probability weighting would obviously not change it!) What does this imply for the objective, i.e. non-weighted, utility that both players obtain? If we consider the common wealth  $U_C$ , i.e. the sum of the utilities of player A and player B,

$$U_C = pq(A_1 + B_1) + p(1 - q)(A_2 + B_2) + (1 - p)(1 - q)(A_3 + B_3),$$

we see that if  $p$  and  $q$  are growing, assumption (i) implies that  $U_C$  is growing as well. In other words, the common wealth is increasing the more the players overweight small probabilities. In a certain sense, overweighting is therefore beneficial for the ‘‘society’’ of players, even though it might not be good from the selfish point of view.

We could also try to explain this result by the following intuitive argument: a possible deviator who overweights the small probability of being ‘‘caught in the act’’ will deviate less, leading to a better common wealth. This simplistic argument, however, is not sound: we could also argue that an enforcer who overweights the small probability of a deviation would check more frequently and hence would cause a lower common wealth, since checking is costly. We see from these fallacious lines of argument that the thorough analysis done up to now was indeed necessary, and can not be replaced with some simple reasoning.

A society, in which all individuals overweight small probabilities, increases its common wealth in social control games (as defined above). In fact, we only need the assumptions (i)-(iii) of definition 2.4 for this result. But what happens if an individual in such a society behaves differently? Does it have an advantage that leads to an erosion of the common overweighting in the population or is probability

overweighting *evolutionarily stable* as in the case of the special game studied in the last section?

In fact, the latter is the case:

**Theorem 2.5.** *Every social control game has the property that every probability weighting  $\gamma \in (0, 1)$  is evolutionarily semi-stable in the sense of Def. 2.1.*

The proof is given in the appendix.

If we consider the dynamics of this meta-game,  $\gamma$  would converge to zero. Again, in any realistic scenario, the players would face more than one type of games which would avoid this difficulty.

### 2.3 Evolutionary stability of probability weighting in games with infinite strategy space

To demonstrate how the analysis of the previous sections can be extended to games with infinitely many pure strategies, we consider the “war of attrition” as first introduced in Bishop & Cannings (1978). Both players decide on a waiting time  $t_i \in [0, \infty)$ . If player A’s waiting time  $t_A$  is longer than his opponent’s then he obtains a utility of  $1 - t_B$ , if it is shorter, he obtains  $-t_A$ . In other words, the players obtain a prize of utility one if they wait longer than their opponents, but have to pay for the time they have to wait. More precisely, the utility of player A waiting  $t_A$  is given by

$$u_A(t_A, t_B) = \begin{cases} 1 - t_B & \text{if } t_A \geq t_B, \\ -t_A & \text{if } t_A < t_B, \\ \frac{1}{2} - t_A & \text{if } t_A = t_B. \end{cases}$$

If both players have rational preferences, i.e. they do not show probability weighting, then the optimal solution is a mixed strategy with probability distribution  $\phi(t) = e^{-t}$  where  $t$  denotes the waiting time (Bishop & Cannings 1978).

To study the influence of probability weighting we compute the prospect utility of player A with probability weighting parameter  $\alpha$  as

$$PT_A(\phi_A, \phi_B) = \frac{1}{\int_0^\infty (\phi_B(s))^\alpha ds} \int_0^\infty \int_0^\infty u_A(t_A, t_B) (\phi_B(t_B))^\alpha dt_B \phi_A(t_A) dt_A,$$

where  $\phi_A$  and  $\phi_B$  denote the mixed strategies of the players<sup>8</sup>.

To compute the Nash equilibrium we need to solve

$$\frac{d}{dt_A} U_A(t_A, \phi_B) = \text{const.},$$

which leads to

$$\phi_B^\alpha(t_A) = -\alpha \phi_B(t_A)^{\alpha-1} \phi_B'(t_A) \alpha.$$

<sup>8</sup>For a derivation of prospect theory for continuous state spaces see Rieger & Wang (2008).

Dividing by  $\phi_B^{\alpha-1}$  and solving the resulting simple differential equation under the side condition that  $\int_0^\infty \phi_B(s) ds = 1$  we obtain

$$\phi_B(t) = \frac{1}{\alpha} e^{-\frac{1}{\alpha}t}.$$

Therefore larger degrees of probability weighting lead on average to shorter waiting times. This reduces average waiting costs for both players (compare Fig. 2) and is therefore beneficial for their common wealth, but from this consideration it is not yet clear, whether probability weighting is also evolutionarily stable or whether a small number of players with a lower degree of probability weighting would outperform a majority of players that have higher probability weighting.

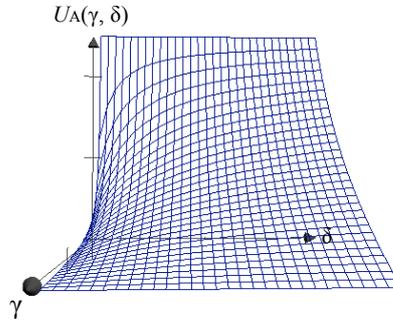


Figure 2: The utility for player A is higher when he has a low  $\gamma$ .

In fact one can prove that this is not the case and probability weighting is indeed evolutionarily stable. As in the case of social control games, one can only prove that probability weighting is semi-evolutionarily stable, i.e. stable against *smaller* degrees of probability weighting in the sense of Def. 2.1 which is sufficient to demonstrate that (irrational) prospect theory preferences can once more survive in the long run. We formulate this result in the following theorem:

**Theorem 2.6.** *Every probability weighting  $\gamma \in (0, 1)$  in the waiting game is semi-evolutionarily stable.*

The proof is given in the appendix. □

### 3 Discussion

Prospect theory describes decisions under risk quite accurately. But how could this be, given that humans are constantly confronted with uncertain situations and rational decisions in such situations should have an evolutionary advantage?

Prospect theory consists basically of two deviations from the rational model of expected utility theory: framing and probability weighting. Framing in gains and losses, rather than in final wealth, seems to be something quite natural if we consider that throughout most of our evolutionary history the accumulation of “wealth” (in whatever sense) was impossible, and it was therefore enough for our ancestors (and it is still enough for animals) to consider only gains and losses. Moreover, it seems very natural for an individual without wealth to be risk-averse in gains (i.e. to assume diminishing marginal utility for gains), but risk-seeking in losses: in most situations, animals cannot afford to “lose”, i.e. to fall behind an average benchmark regarding food or health, since this poses in a competitive environment a high risk of dying, so that it makes sense to try to avoid any loss even by taking high risks. This idea can explain why humans and animals frame in gains and losses (McDermott, Fowler & Smirnov 2008). The puzzle of probability weighting in humans, however, remains.

There have been approaches to answer this problem, most notably from the viewpoint of psychophysics, compare Tversky & Kahneman (1992). They argue that decisions on low probability events are more frequent and that therefore differences between such low probabilities are overweighted. However, this does not explain the evolutionary advantage such an overweighting should give.

But maybe this “probability puzzle” can be resolved with a very different observation, namely that probability weighting can have positive effects if individuals have interactions with each other. We have seen that this is the case when we have social interaction games or games like the war of attrition: here individuals of a society *profit* from probability weighting of their members. Probability weighting becomes evolutionarily stable.

The ideas presented here are of course a somehow speculative suggestion how to explain the probability puzzle.

One could question this approach for explaining the probability weighting puzzle by arguing that humans face many different games and different decision situations in their lives, not only social control games or the war of attrition. In some of them overweighting of small probabilities might be evolutionarily advantageous, in others it might be disadvantageous. Of course, we cannot prove which of these situations is more important. However, noticing that there are many situations in which probability weighting in *some direction* is useful, makes it unlikely that “on average” both situations exactly balance and a neutral weighting should be optimal. On the contrary, it seems natural that *some* deviation from a neutral weighting will be observed. Experiments tell us that this deviation tends to go into the direction that overweights small probabilities (as it would be optimal, e.g., in the case of social control games or the war of attrition). The assumption that neutral weighting is evolutionarily optimal is *only* natural as long as we neglect interactions between individuals. As soon as we take problems in game theory into account, probability weighting becomes in the generic case optimal.

## A Proofs

PROOF OF THEOREM 2.5:

Let  $\delta > \gamma$ . We verify the inequality (6) by considering the difference of both sides:

$$\begin{aligned}\Delta_\varepsilon &:= \varepsilon U(\gamma, \delta) + (1 - \varepsilon)U(\gamma, \gamma) - \varepsilon U(\delta, \delta) - (1 - \varepsilon)U(\delta, \gamma) \\ &= \frac{1}{2} \left( \varepsilon(U_A(\gamma, \delta) + U_B(\delta, \gamma)) + (1 - \varepsilon)(U_A(\gamma, \gamma) + U_B(\gamma, \gamma)) \right. \\ &\quad \left. - \varepsilon(U_A(\delta, \delta) + U_B(\delta, \delta)) - (1 - \varepsilon)(U_A(\delta, \gamma) + U_B(\gamma, \delta)) \right).\end{aligned}$$

We denote the PT Nash equilibrium strategies of the players by  $p_\beta$  and  $q_\alpha$  if player A is an  $\alpha$ -weighter and player B is a  $\beta$ -weighter. With this we can write

$$U_A(\gamma, \delta) = p_\delta q_\gamma A_1 + p_\delta(1 - q_\gamma)A_3 + (1 - p_\delta)q_\gamma A_2 \quad \text{etc.}$$

We prove that  $\Delta_0 > 0$ . After a small calculation we arrive at

$$\begin{aligned}\Delta_0 &= \frac{1}{2} \left( (A_1 - A_2 - A_3)(q_\gamma - q_\delta)p_\gamma + (B_1 - B_2 - B_3)(p_\gamma - p_\delta)q_\gamma \right. \\ &\quad \left. + A_2(q_\gamma - q_\delta) + B_3(p_\gamma - p_\delta) \right).\end{aligned}\tag{7}$$

We obtain from the definition of a social control game that  $A_1 - A_2 - A_3 > 0$ ,  $B_1 - B_2 - B_3 < 0$ , moreover we have already seen that  $p_\gamma - p_\delta > 0$ ,  $q_\gamma - q_\delta > 0$  and  $q_\gamma < 1$ . Hence we can estimate (7) as

$$\Delta_0 \geq \frac{1}{2} \left( A_2(q_\gamma - q_\delta) + (B_1 - B_2)(p_\gamma - p_\delta) \right).$$

using the initial assumptions (i) and (ii), we see that  $\Delta_0 > 0$ . Now, since  $\Delta_\varepsilon$  is continuous, we deduce that, for  $\varepsilon > 0$  sufficiently small,  $\Delta_\varepsilon > 0$ . This proves inequality (6).  $\square$

PROOF OF THEOREM 2.6:

The expected utility of player A with probability weighting  $\alpha$  when playing against a player with probability weighting  $\beta$  is given by

$$U_A(\alpha, \beta) = \int_0^\infty \int_0^{t_A} (1 - t_B) \frac{1}{\alpha} e^{-\frac{1}{\alpha} t_B} dt_B + \int_{t_A}^\infty -t_A \frac{1}{\alpha} e^{-\frac{1}{\alpha} t_B} dt_B \frac{1}{\beta} e^{-\frac{1}{\beta} t_A} dt_A.$$

A small computation gives

$$\begin{aligned}U_A(\alpha, \beta) &= \int_0^\infty \left( 1 - e^{-\frac{1}{\alpha} t_A} + \alpha e^{-\frac{1}{\alpha} t_A} - \alpha \right) \frac{1}{\beta} e^{-\frac{1}{\beta} t_A} dt_A \\ &= 1 - \alpha + \frac{\alpha - 1}{\left( \frac{1}{\alpha} + \frac{1}{\beta} \right) \beta} \\ &= \frac{\beta - \alpha \beta}{\alpha + \beta}.\end{aligned}$$

Now let  $0 < \gamma < \delta \leq 1$  and  $\varepsilon > 0$ , then

$$\begin{aligned} (1 - \varepsilon)U_A(\gamma, \gamma) + \varepsilon U_A(\gamma, \delta) &= (1 - \varepsilon)\frac{1 - \gamma}{2} + \varepsilon\frac{\delta - \gamma\delta}{\gamma + \delta} \\ &> (1 - \varepsilon)\frac{\gamma - \gamma\delta}{\gamma + \delta} + \varepsilon\frac{\delta - \delta^2}{2\delta} = (1 - \varepsilon)U_A(\delta, \gamma) + \varepsilon U_A(\delta, \delta). \end{aligned}$$

Therefore  $\gamma$  is semi-evolutionarily stable.  $\square$

## Acknowledgement

I thank Anke Gerber and Frank Riedel for their very valuable suggestions regarding this work, Lars Koch for his helpful comments, Mei Wang for interesting discussions which initiated this work and Thorsten Hens for his steady support. The support by the National Centre of Competence in Research "Financial Valuation and Risk Management" (NCCR FINRISK), Project 3, "Evolution and Foundations of Financial Markets", and by the University Priority Program "Finance and Financial Markets" of the University of Zürich is gratefully acknowledged.

## References

- Bishop, D. & Cannings, C. (1978), 'A generalized war of attrition', *Journal of Theoretical Biology* **70**(1), 85–124.
- Dekel, E., Safra, Z. & Segal, U. (1991), 'Existence and dynamic consistency of Nash equilibrium with non-expected utility preferences', *Journal of Economic Theory* **55**, 229–246.
- Ely, J. C. & Yilankaya, O. (2001), 'Nash equilibrium and the evolution of preferences', *Journal of Economic Theory* **97**(2), 255–272.
- Gilboa, I. & Schmeidler, D. (1992), Updating ambiguous beliefs, in 'TARK '92: Proceedings of the 4th conference on Theoretical aspects of reasoning about knowledge', Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, pp. 143–162.
- Goeree, J. K., Holt, C. A. & Palfrey, T. R. (2003), 'Risk averse behavior in generalized matching pennies games', *Games and Economic Behavior* **45**(1), 97–113.
- Güth, W. & Yaari, M. E. (1992), Explaining reciprocal behavior in simple strategic games: an evolutionary approach, in U. Witt, ed., 'Explaining Process and Change: Approaches to Evolutionary Economics', University of Michigan Press, pp. 23–34.
- Kahneman, D. & Tversky, A. (1979), 'Prospect Theory: An analysis of decision under risk', *Econometrica* **47**, 263–291.
- McDermott, R., Fowler, J. H. & Smirnov, O. (2008), 'On the evolutionary origin of prospect theory preferences', *Journal of Politics*.
- Rieger, M. O. & Koch, L. (2009), Prospect theory preferences in games, Work in progress.
- Rieger, M. O. & Wang, M. (2008), 'Prospect Theory for continuous distributions', *Journal of Risk and Uncertainty* **36**, 83–102.

- Schmeidler, D. (1989), 'Subjective probability and expected utility without additivity', *Econometrica* **57**(3), 571–587.
- Schneider, S. L. & Lopes, L. L. (1986), 'Reflection in preferences under risk: who and when may suggest why', *Journal of Experimental Psychology: Human Perception and Performance* **12**, 535–548.
- Tversky, A. & Kahneman, D. (1992), 'Advances in Prospect Theory: Cumulative representation of uncertainty', *Journal of Risk and Uncertainty* **5**, 297–323.
- von Neumann, J. (1928), 'Zur Theorie der Gesellschaftsspiele', *Mathematische Annalen* **100**(1), 295–320.
- von Neumann, J. & Morgenstern, O. (1944), *Theory of Games and Economic Behavior*, Princeton University Press, Princeton, NJ.
- Wakker, P. P. (1989), Transforming probabilities without violating stochastic dominance, in 'Mathematical Psychology in Progress', Springer, Berlin, pp. 29–47.