



HAL
open science

REVERSIBILITY AND OSCILLATIONS IN ZERO-SUM DISCOUNTED STOCHASTIC GAMES

Sylvain Sorin, Guillaume Vigeral

► **To cite this version:**

Sylvain Sorin, Guillaume Vigeral. REVERSIBILITY AND OSCILLATIONS IN ZERO-SUM DISCOUNTED STOCHASTIC GAMES. 2013. hal-00869656v1

HAL Id: hal-00869656

<https://hal.science/hal-00869656v1>

Preprint submitted on 4 Oct 2013 (v1), last revised 31 Oct 2013 (v2)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

REVERSIBILITY AND OSCILLATIONS IN ZERO-SUM DISCOUNTED STOCHASTIC GAMES

SYLVAIN SORIN AND GUILLAUME VIGERAL

ABSTRACT. We show that by coupling two well-behaved exit-time problems one can construct two-person zero-sum stochastic games with finite state space having oscillating discounted values. This unifies and generalizes recent examples due to Vigerál (2013) and Ziliotto (2013).

CONTENTS

1. Introduction	2
2. A basic model	2
2.1. 0 player	2
2.2. 1 player	3
2.3. 2 players	3
3. Reversibility	4
3.1. Two examples	4
3.2. The discounted framework	5
4. Some regular configurations of order $\frac{1}{2}$	6
4.1. A regular configuration with 0 players and countable state space	6
4.2. A regular configuration with one player, finitely many states, compact action space and continuous transition	7
4.3. A regular configuration with two players and finitely many states and actions	7
5. Some oscillating configurations of order $\frac{1}{2}$	7
5.1. Example 4.2. perturbed	7
5.2. Example 4.3. perturbed	8
5.3. Countable action space	8
5.4. Countable state space	9
5.5. A MDP with signals	10
5.6. A game in the dark	10
6. Combinatorics	10
6.1. Example 4.2 + Example 5.1	10
6.2. Examples 4.2 + Example 5.3	10
6.3. Example 4.3 + Example 5.1	11
6.4. Example 5.2 + 5.2	11
6.5. Example 5.4 + 5.4, 5.5 + 5.5 and 5.6 + 5.6	11
6.6. Example 4.1 + 5.4	11
7. Comparison and conclusion	11
7.1. Irreversibility	11
7.2.	11
7.3. Semi-algebraic	11
7.4. Related issues	11
References	12

1. INTRODUCTION

1) We construct a family of zero-sum games in discrete time where the purpose is to control the law of a stopping time of exit. For each evaluation of the stream of outcomes (defined by a probability distribution on the positive integers $n = 1, 2, \dots$), value and optimal strategies are well defined.

In particular for a given discount factor $\lambda \in]0, 1]$ optimal stationary strategies define an inertia rate Q_λ .

When two such configurations (1 and 2) are coupled this induces a stochastic game where the state will move from one to the other in a way depending on the previous rates $Q_\lambda^i, i = 1, 2$.

The main observation is that the discounted value is a function of the ratio $\frac{Q_\lambda^1}{Q_\lambda^2}$ that can oscillate as λ goes to 0, when both inertia rates converge to 0.

2) This construction reveals a common structure in two recent “counter-examples” by Vigeral [12] and Ziliotto [13] dealing with two-person zero-sum stochastic games with finite state space: compact action spaces and standard signalling in the first case, finite action spaces and signals on the state space in the second. In both cases it was proved that the family of discounted values does not converge.

2. A BASIC MODEL

A *configuration* P is defined by a general two-person repeated game in discrete time (see [7]) on a state space Ω with a specific starting state $\bar{\omega}$ and a subset $\bar{\Omega}$ satisfying $\bar{\omega} \in \bar{\Omega} \subset \Omega$.

Let S be the stopping time of exit of $\bar{\Omega}$:

$$S = \min\{n \in \mathbb{N}; \omega_n \notin \bar{\Omega}\}$$

where ω_n is the state at stage n .

Each couple of strategies (σ, τ) of the players specifies, with the parameters of the game (initial state, transition function on states and signals), the law of S . For each evaluation $\theta = \{\theta_n\}$ on the set of positive integers $\mathbb{N}^* = 1, 2, \dots$, let $d_\theta(\sigma, \tau)$ be the expected (normalized) duration spent in $\bar{\Omega}$:

$$d_\theta(\sigma, \tau) = E_{\sigma, \tau} \left[\sum_{n=1}^{S-1} \theta_n \right].$$

For each real parameters $\alpha < \beta$, consider the game with payoff α at any state in $\bar{\Omega}$ and with absorbing payoff β in its complement $\Omega \setminus \bar{\Omega}$.

Then for any evaluation θ , Player 1 (the maximizer) minimizes $d_\theta(\sigma, \tau)$ since the payoff $\gamma_\theta(\sigma, \tau)$ is given by:

$$\gamma_\theta(\sigma, \tau) = \alpha d_\theta(\sigma, \tau) + \beta(1 - d_\theta(\sigma, \tau)).$$

Lemma 2.1. *In particular if the game has a value v_θ then*

$$v_\theta = \alpha Q_\theta + \beta(1 - Q_\theta)$$

with $Q_\theta = \sup_\tau \inf_\sigma g_\theta(\sigma, \tau) = \inf_\sigma \sup_\tau g_\theta(\sigma, \tau)$, called the inertia rate.

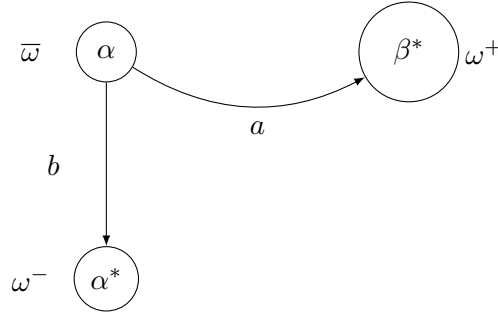
Here are 3 examples corresponding to a Markov Chain (0 player), a Dynamic Programming Problem (1 player) and a Stochastic Game (2 players).

In all cases $\Omega = \{\bar{\omega}, \omega^+, \omega^-\}$ and $\bar{\Omega} = \{\bar{\omega}, \omega^-\}$, hence S is the first time where the exit state ω^+ is reached. Moreover ω^- is an absorbing state.

2.1. 0 player.

a resp. b is the probability to go from $\bar{\omega}$ to ω^+ (resp. to ω^-) with $a, b, a + b \in [0, 1]$.

FIGURE 1

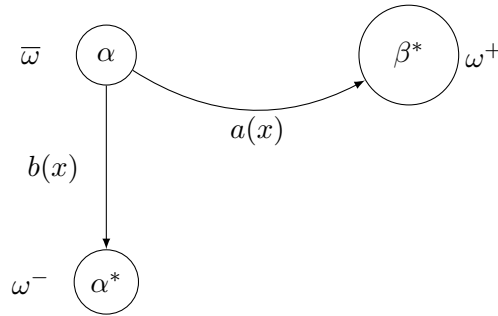


a “*” stands for an absorbing payoff.

2.2. 1 player.

The action set is $X = [0, 1]$ and the impact of an action x is on the transitions, given by $a(x)$ from $\bar{\omega}$ to ω^+ and $b(x)$ from $\bar{\omega}$ to ω^- , where a and b are two continuous function from $[0, 1]$ to $[0, 1]$ with $a + b \in [0, 1]$.

FIGURE 2



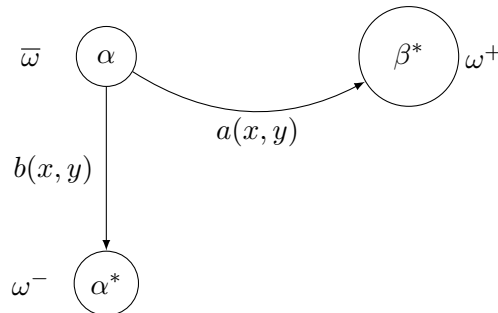
2.3. 2 players.

In state $\bar{\omega}$ the players have two actions and the transitions are given by:

	Stay	Quit
Stay	$\bar{\omega}$	ω^+
Quit	ω^+	ω^-

Let x (resp. y) be the probability on Stay and $a(x, y) = x(1 - y) + y(1 - x)$, $b(x, y) = xy$. The mixed extension gives the *configuration*:

FIGURE 3



Of course one can define such a configuration for any maps a and b from $[0, 1]^2$ to $[0, 1]$.

Consider the λ -discounted case P_λ . Let $r_\lambda(x, y)$ be the induced expected payoff.

Lemma 2.2.

$$r_\lambda(x, y) = \frac{(\lambda + (1 - \lambda)b(x, y)) \times \alpha + (1 - \lambda)a(x, y) \times \beta}{\lambda + (1 - \lambda)(a(x, y) + b(x, y))}.$$

Proof.

By stationarity:

$$r_\lambda(x, y) = \lambda \times \alpha + (1 - \lambda)[a(x, y) \times \beta + b(x, y) \times \alpha + (1 - a(x, y) - b(x, y)) \times r_\lambda(x, y)]$$

□

In particular letting:

$$(1) \quad q_\lambda(x, y) = \frac{(\lambda + (1 - \lambda)b(x, y))}{\lambda + (1 - \lambda)(a(x, y) + b(x, y))}$$

one has:

$$(2) \quad r_\lambda(x, y) = q_\lambda(x, y) \times \alpha + (1 - q_\lambda(x, y)) \times \beta.$$

In the normalized game of length one, $q_\lambda(x, y)$ is the expected duration spent with payoff α before reaching the absorbing state $\bar{\omega}$ with payoff β .

Lemma 2.3.

$Q_\lambda = \min_x \max_y q_\lambda(x, y) = \max_y \min_x q_\lambda(x, y)$ and the value v_λ of P_λ satisfies:

$$(3) \quad v_\lambda = Q_\lambda \times \alpha + (1 - Q_\lambda) \times \beta.$$

Note that the value exists either in the one player case or when a and b are bilinear (and hence q_λ is quasi concave/convex).

3. REVERSIBILITY

Consider now a two person zero-sum stochastic game G generated by two dual configurations P^1 and P^2 of the previous type, with $\alpha^1 = -1$ and $\alpha^2 = 1$, which are coupled in the following sense: the exit domain from P^1 ($\Omega^1 \setminus \bar{\Omega}^1$) is the starting state $\bar{\omega}^2$ in P^2 and reciprocally. In addition we assume that the exit events are known by both players and that both configurations have a value.

3.1. Two examples.

3.1.1. Two configurations with one player in each.

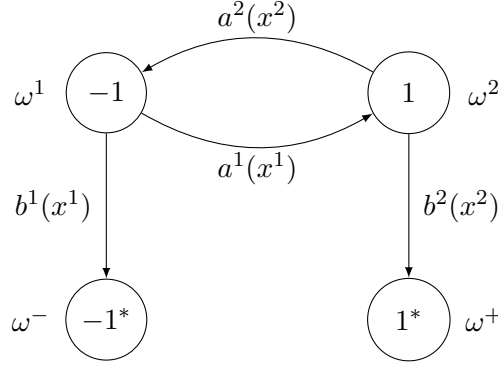
There are four states $\Omega = \{\omega^1, \omega^2, \omega^-, \omega^+\}$.

Both ω^+ and ω^- are absorbing states with constant payoff $+1$ and -1 , respectively.

The payoff in state ω^i is also constant and equals to -1 for $i = 1$ and to $+1$ for $i = 2$. The action set for player 1 is $X^1 = [0, 1]$ and the impact of an action x^1 on the transitions is given by $a^1(x^1)$ from ω^1 to ω^2 and $b^1(x^1)$ from ω^1 to ω^- , where a^1 and b^1 are two continuous function from $[0, 1]$ to $[0, 1]$.

Similarly the action set for player 2 is $X^2 = [0, 1]$ and $a^2(x^2)$ is the transition probability from ω^2 to ω^1 and $b^2(x^2)$ from ω^2 to ω^+ .

FIGURE 4



3.1.2. Two configurations with 2 players.

There are two absorbing states with payoff 1 and -1 . In the two other states (ω^1 and ω^2) the payoff is constant and the transitions are given by the following matrices (compare to Bewley and Kohlberg [1]):

ω^2	Stay	Quit
Stay	1	$\xrightarrow{1}$
Quit	$\xrightarrow{1}$	1^*

ω^1	Stay	Quit
Stay	-1	$\xleftarrow{-1}$
Quit	$\xleftarrow{-1}$	-1^*

where an arrow means a transition to the other state.

3.2. The discounted framework.

For each $\lambda \in]0, 1[$ the coupling between the two configurations defines a discounted game G_λ with value v_λ satisfying:

$$v_\lambda(\bar{\omega}^1) = v_\lambda^1 \in [-1, 1[, \quad v_\lambda(\bar{\omega}^2) = v_\lambda^2 \in]-1, 1].$$

In particular, starting from state $\bar{\omega}^1$ the model is equivalent to the one with an exit state $\bar{\omega}^2$ with absorbing payoff $v_\lambda(\bar{\omega}^2)$ (by stationarity of the evaluation), which thus corresponds to the payoff $\beta_1 > \alpha_1$ in the configuration P of the previous section 2.

Hence one obtains, using Lemma 2.1, that $\{v_\lambda^i\}$ is a solution of the next system of equations:

Proposition 3.1.

$$\begin{aligned} v_\lambda^1 &= Q_\lambda^1 \times (-1) + (1 - Q_\lambda^1) \times v_\lambda^2 \\ v_\lambda^2 &= Q_\lambda^2 \times (+1) + (1 - Q_\lambda^2) \times v_\lambda^1. \end{aligned}$$

It follows that:

Corollary 3.1.

$$\begin{aligned} v_\lambda^1 &= \frac{Q_\lambda^2 - Q_\lambda^1 - Q_\lambda^1 Q_\lambda^2}{Q_\lambda^1 + Q_\lambda^2 - Q_\lambda^1 Q_\lambda^2} \\ v_\lambda^2 &= \frac{Q_\lambda^2 - Q_\lambda^1 + Q_\lambda^1 Q_\lambda^2}{Q_\lambda^1 + Q_\lambda^2 - Q_\lambda^1 Q_\lambda^2} \end{aligned}$$

Comments:

- 1) As λ goes to 0, Q_λ converges to 0 in the model of section 2.2, as soon as $\limsup \frac{a(x)}{b(x)} = +\infty$, as x goes to 0.

2) In the current framework, assuming that both Q_λ^i go to 0, the asymptotic behavior of v_λ^1 depends upon the evolution of the ratio $\frac{Q_\lambda^1}{Q_\lambda^2}$. In fact one has:

$$v_\lambda^1 \sim v_\lambda^2 \sim \frac{1 - \frac{Q_\lambda^1}{Q_\lambda^2}}{1 + \frac{Q_\lambda^1}{Q_\lambda^2}}.$$

3) In particular one obtains:

Theorem 3.1. *Assume that both Q_λ^i go to 0 as λ goes to 0 and that $\frac{Q_\lambda^1}{Q_\lambda^2}$ has more than one accumulation point, then v_λ^i does not converge.*

More precisely it is enough that $Q_\lambda^i \sim \lambda^r f^i(\lambda)$ for some $r > 0$, with $0 < A \leq f^i \leq B$ and that one of the $f^i(\lambda)$ does not converge as λ goes to 0, to obtain the result.

The next section 4 will describe several models generating such probabilities Q_λ^i , with f^i converging or not.

We will use the terminology *regular or oscillating configurations*.

The above result implies that by coupling any two of these configurations (of the same order of magnitude r) where one is oscillating, one can generate a stochastic game for which the family of discounted values does not converge, see Section 6. In the next two sections we give examples of, respectively, regular or oscillating configurations of order $\frac{1}{2}$.

4. SOME REGULAR CONFIGURATIONS OF ORDER $\frac{1}{2}$

We give here three examples of regular configurations of order $\frac{1}{2}$. Let us remark right now that these configurations are in a certain sense minimal ones. Any configuration with one player, with finitely many states and actions and full observation is, by Blackwell optimality, asymptotically equivalent to a finite Markov chain. And in any such chain,

- either with positive probability there is no exit, and Q_λ is of order 0.
- or at each stage, given no prior exit there is exit in the next m stages with probability at least p , where m and $p > 0$ are fixed. This implies that Q_λ is of order 1.

4.1. A regular configuration with 0 players and countable state space.

Consider a random walk on $\mathbf{N} \cup \{-1\}$ and exit state -1 . In any other state $m \in \mathbf{N}$ the transition is $\frac{1}{2}\delta_{m-1} + \frac{1}{2}\delta_{m+1}$. The starting state is 0. Denote by s_n the probability that exit happens at stage n ; it is well known (theorem 5b p 164 in [3]) that the generating function of S is given by $F(z) = \frac{1 - \sqrt{1 - z^2}}{z}$. Hence,

$$\begin{aligned} Q_\lambda &= \sum_{n=1}^{+\infty} s_n \sum_{t=1}^n \lambda(1-\lambda)^{i-1} \\ &= \sum_{n=1}^{+\infty} s_n (1 - (1-\lambda)^n) \\ &= F(1) - F(1-\lambda) \\ &= \frac{\sqrt{2\lambda - \lambda^2} - \lambda}{1-\lambda} \\ &\sim \sqrt{2\lambda}. \end{aligned}$$

4.2. A regular configuration with one player, finitely many states, compact action space and continuous transition.

Consider example 2.2.

Take $a(x) = x$ and $b(x) = x^2$. Then $Q_\lambda = \min_x \left\{ \frac{\lambda + (1-\lambda)x^2}{\lambda + (1-\lambda)x^2 + (1-\lambda)x} \right\}$ and a first order condition gives $x_\lambda = \sqrt{\frac{\lambda}{1-\lambda}}$ hence $Q_\lambda \sim 2\sqrt{\lambda}$.

4.3. A regular configuration with two players and finitely many states and actions.

Consider example 2.3.

It is straightforward [12] to compute that in Γ_λ the optimal strategy for each player is $x_\lambda = y_\lambda = \frac{\sqrt{\lambda}}{1+\sqrt{\lambda}}$. Hence:

$$\begin{aligned} Q_\lambda &= \frac{\lambda + (1-\lambda)x_\lambda y_\lambda}{\lambda + (1-\lambda)(x_\lambda + y_\lambda - x_\lambda y_\lambda)} \\ &\sim \sqrt{\lambda}. \end{aligned}$$

5. SOME OSCILLATING CONFIGURATIONS OF ORDER $\frac{1}{2}$

5.1. Example 4.2. perturbed.

Recall that the choice of $a(x) = x$ and $b(x) = x^2$ leads to $Q_\lambda \sim 2\sqrt{\lambda}$.

To get oscillations one can choose $b = x^2$ and $a(x) = xf(x)$ with $f(x)$ bounded away from 0, oscillating and such that $f'(x) = o(1/x)$. For example, $f(x) = 2 + \sin(\ln(-\ln x))$.

Proposition 5.1. *For this choice of transition functions one has:*

$$Q_\lambda \sim \frac{2\sqrt{\lambda}}{f(\sqrt{\lambda})}.$$

Proof.

In fact recall by (1) that x_λ minimizes $q_\lambda(x)$ iff it minimizes $\rho_\lambda(x) = \frac{\lambda + (1-\lambda)b(x)}{a(x)}$ and then $Q_\lambda \sim \rho_\lambda(x_\lambda)$ as soon as they both tend to 0.

The first order condition gives:

$$\frac{\lambda}{1-\lambda} = \frac{x^2(f(x) - xf'(x))}{f(x) + xf'(x)}$$

which leads to:

$$x_\lambda \sim \sqrt{\lambda}.$$

By the mean value theorem and since $f'(x) = o(1/x)$,

$$\frac{\|f(x_\lambda) - f(\sqrt{\lambda})\|}{\|x_\lambda - \sqrt{\lambda}\|} = o\left(\frac{1}{\sqrt{\lambda}}\right)$$

hence $f(x_\lambda) \sim f(\sqrt{\lambda})$ and:

$$Q_\lambda \sim \frac{2\sqrt{\lambda}}{f(\sqrt{\lambda})}.$$

In particular $\frac{Q_\lambda}{\sqrt{\lambda}}$ has not limit. □

5.2. Example 4.3. perturbed.

Let $s \in C^1([0, \frac{1}{16}], \mathbb{R})$ such that s and $x \rightarrow xs'(x)$ are both bounded by $\frac{1}{16}$. Consider a configuration as in FIGURE 3 but for perturbed functions a and b :

$$\begin{aligned} a(x, y) &= \frac{(\sqrt{x} + \sqrt{y})(1 - \sqrt{x} + s(x))(1 - \sqrt{y} + s(y))}{2(1-x)(1-y)(1-f_2(x, y))} \\ b(x, y) &= \frac{\sqrt{xy} [(1 - \sqrt{x})(1 - \sqrt{y}) + f_1(x, y) - \sqrt{xy}f_2(x, y)]}{(1-x)(1-y)(1-f_2(x, y))}. \end{aligned}$$

where

$$f_1(x, y) = \begin{cases} \frac{\sqrt{x}s(x) - \sqrt{y}s(y)}{\sqrt{x} - \sqrt{y}} & \text{if } x \neq y \\ 2xs'(x) + s(x) & \text{if } x = y \end{cases}$$

and

$$f_2(x, y) = \begin{cases} \frac{\sqrt{y}s(x) - \sqrt{x}s(y)}{\sqrt{x} - \sqrt{y}} & \text{if } x \neq y \\ 2xs'(x) - s(x) & \text{if } x = y \end{cases}$$

Then a and b are continuous (Lemma 12 and Lemma 10 in [12]) and $x_\lambda = y_\lambda = \lambda$ are optimal in the game with payoff q_λ [12]. Hence:

$$\begin{aligned} Q_\lambda &= \frac{\frac{\lambda}{1-\lambda} + b(\lambda, \lambda)}{\frac{\lambda}{1-\lambda} + b(\lambda, \lambda) + a(\lambda, \lambda)} \\ &\sim \frac{\lambda + \frac{\lambda(1+s(\lambda)+2\lambda s'(\lambda))}{1+s(\lambda)-2\lambda s'(\lambda)}}{\frac{2\sqrt{\lambda}(1+s(\lambda))^2}{2(1+s(\lambda)-2\lambda s'(\lambda))}} \\ &\sim \frac{2\sqrt{\lambda}}{1+s(\lambda)} \\ &\sim \frac{\lambda(1+s(\lambda) - 2\lambda s'(\lambda) + 1 + s(\lambda) + 2\lambda s'(\lambda))}{\sqrt{\lambda}(1+s(\lambda))^2} \end{aligned}$$

The configuration is thus oscillating for $s(x) = \frac{\sin \ln x}{16}$ for example.

Next we recall 4 models that appears in Ziliotto [13] (in which the divergence of v_λ was proven) and we compute the corresponding Q_λ .

5.3. Countable action space.

Consider again the example 2.2 but assume now that the action space X is \mathbb{N}^* and no longer $[0, 1]$. The transition are given by $(a_n, b_n) = (\frac{1}{2^n}, \frac{1}{4^n})$.

Proposition 5.2.

For this configuration $Q_\lambda/\sqrt{\lambda}$ oscillates on a sequence $\{\lambda_m\}$ of discount factors like $\lambda_m = \frac{1}{2^m}$.

Proof.

Note first that the choice n inducing $x = \frac{1}{2^n}$ is asymptotically optimal for $a(x) = x, b(x) = x^2$, like in example 4.2, at $\lambda = \frac{1}{4^n}$ and $Q_\lambda \sim 2\sqrt{\lambda}$.

For $\lambda^2 = \frac{1}{4^n} \frac{1}{4^{n+1}}$ one obtains:

$$\begin{aligned} \rho_\lambda\left(\frac{1}{2^n}\right) &\sim \left(\frac{1}{2} \times \frac{1}{4^n} + \frac{1}{4^n}\right)2^n \\ &\sim \frac{3\sqrt{2}}{2}\sqrt{\lambda} \end{aligned}$$

and similarly:

$$\begin{aligned}\rho_\lambda\left(\frac{1}{2^{n+1}}\right) &\sim \left(\frac{1}{2} \times \frac{1}{4^n} + \frac{1}{4^{n+1}}\right)2^{n+1} \\ &\sim \frac{3\sqrt{2}}{2}\sqrt{\lambda}.\end{aligned}$$

Finally one checks that $\rho_\lambda\left(\frac{1}{2^{n+m}}\right) \geq \frac{3\sqrt{2}}{2}\sqrt{\lambda}$ for $m = -n, \dots, -1$ and $m \geq 2$.

Thus for this specific λ , $\rho_\lambda(x)$ is bounded below by a quantity of the order $\frac{3\sqrt{2}}{2}\sqrt{\lambda}$.

It follows that $Q_\lambda/\sqrt{\lambda}$ oscillates between 2 and $\frac{3\sqrt{2}}{2}$ on a sequence $\{\lambda_m\}$ of discount factors like $\lambda_m = \frac{1}{2^m}$. \square

Note that this result is conceptually similar to example 5.1.

5.4. Countable state space.

We consider here a configuration which is the dual of the previous one with now finite action space and countably many states.

The state space is a countable family of probabilities $y = (y^A, y^B)$ on two positions A and B with $y_n = (\frac{1}{2^n}, 1 - \frac{1}{2^n})$, $n = 0, 1, \dots$, and two absorbing states A^* and B^* .

The player has two actions: *Stay* or *Quit*. Consider state y_n . Under *Quit* an absorbing state is reached: A^* with probability y_n^A and B^* with probability y_n^B . Under *Stay* the state evolves from y_n to y_{n+1} with probability $1/2$ and to $y_0 = (1, 0)$ with probability $1/2$.

The player is informed upon the state, a the starting state is y_0 and the exit state is B^* .

A strategy of the player can be identified with a stopping time corresponding to the first state y_n when he chooses *Quit*.

Let T_n be the random time corresponding to the first occurrence of y_n (under *Stay*) and μ_n the associated strategy: *Quit* (for the first time) at y_n .

Proposition 5.3.

Under μ_n the λ -discounted normalized duration before B^* is

$$q_\lambda(n) = 1 - \frac{(1 - \lambda^2) \left(1 - \frac{1}{2^n}\right)}{1 + 2^{n+1}\lambda(1 - \lambda)^{-n} - \lambda}$$

Proof.

Lemma 2.5 in Ziliotto [13] gives

$$\mathbb{E}[(1 - \lambda)^{T_n}] = \frac{1 - \lambda^2}{1 + 2^{n+1}\lambda(1 - \lambda)^{-n} - \lambda}$$

and

$$q_\lambda(n) = 1 + \left(\frac{1}{2^n} - 1\right)\mathbb{E}[(1 - \lambda)^{T_n}].$$

\square

Proposition 5.4. *The configuration is irregular : $\frac{Q_\lambda}{\sqrt{\lambda}}$ oscillates between two positive values.*

Proof.

With our notations, Ziliotto's Lemma 2.8 [13] states that

$$q_\lambda\left(-\frac{\ln \lambda + \ln 2 + 2 \ln c}{2 \ln 2}\right) \sim (c + c^{-1})\sqrt{2\lambda}.$$

Hence asymptotically,

$$\frac{Q_\lambda}{\sqrt{2\lambda}} = \min \left\{ c + c^{-1} \mid -\frac{\ln \lambda + \ln 2 + 2 \ln c}{2 \ln 2} \in \mathbf{N} \right\}$$

When $-\frac{\ln \lambda + \ln 2}{2 \ln 2}$ is an integer, one can take $c = 1$ which gives $Q_\lambda \sim 2\sqrt{2\lambda}$. Whereas when $-\frac{\ln \lambda + \ln 2}{2 \ln 2}$ is an integer plus one half, the best choice is $c = \sqrt{2}$, leading to $Q_\lambda \sim 3\sqrt{\lambda}$ \square

5.5. A MDP with signals.

The next configuration corresponds to a Markov decision process with 2 states: A and B , 2 absorbing states A^* and B^* and with signals on the state. The player has 2 actions: *Stay* or *Quit*. The transition are as follows:

	A	$\frac{1}{2}; \ell$	$\frac{1}{2}; r$		B	$\frac{1}{2}; \ell$	$\frac{1}{2}; r$
Stay	A	$(\frac{1}{2}A + \frac{1}{2}B)$		Stay	A	B	
Quit	A^*	A^*		Quit	B^*	B^*	

Hence the transition is random: with probability $1/2$ of type ℓ and probability $1/2$ of type r . The player is not informed upon the state reached but only on the signal ℓ or r .

The natural “auxiliary state” space is then the beliefs of the player on (A, B) and one can check [13] that the model is equivalent to the previous one, starting from A and where the exit state is B^* . In fact under Stay, ℓ occurs with probability $1/2$ and the new parameter is $y_0 = (1, 0)$. On the other hand, after r the belief evolves from y_n to y_{n+1} .

Again this configuration generates an oscillating Q_λ of the order of $\sqrt{\lambda}$.

5.6. A game in the dark.

A next transformation is to introduce two players and to generate the random variable $\frac{1}{2}(\ell) + \frac{1}{2}(r)$ in the above model by a process induced by the moves of the players.

This leads to the original framework of the game defined by Ziliotto [13]: action and state spaces are finite and the only information of the players is the initial state and the sequence of moves along the play.

Player 1 has three moves: Stay1, Stay2 and Quit, and player 2 has 2 moves: Left and Right. The payoff is -1 and the transition are as follows:

	A	$Left$	$Right$		B	$Left$	$Right$
Stay1	A	$(\frac{1}{2}A + \frac{1}{2}B)$		Stay1	A	B	
Stay2	$(\frac{1}{2}A + \frac{1}{2}B)$	A		Stay2	B	A	
Quit	A^*	A^*		Quit	B^*	B^*	

By playing $(1/2, 1/2, 0)$ (resp. $(1/2, 1/2)$) player 1 (resp. player 2) can mimick the previous distribution on (ℓ, r) where ℓ corresponds to the event “the moves are on the main diagonal”. It follows that this behavior is consistent with optimal strategies hence the induced distribution on plays is like in the previous example 5.5.

6. COMBINATORICS

In order to obtain oscillations for the discounted values of a stochastic game, it is enough to consider the coupled dynamics generated by a regular and an oscillating configuration, both of order $\frac{1}{2}$.

6.1. Example 4.2 + Example 5.1.

Combining these two configurations yields a coupling of two one-person decision problems, hence a compact stochastic game with perfect information and no asymptotic value. Remark that the transition functions can be taken as smooth as one wants.

6.2. Examples 4.2 + Example 5.3. With this combination one recovers exactly an example of Ziliotto (see section 4.2 in [13]) which is also a stochastic game with perfect information and no asymptotic value. The main difference is that in that case the action space of Player 1 is countable instead of being an interval.

6.3. Example 4.3 + Example 5.1.

Combining these two configurations yields a stochastic game with finite action space for player 2 and no asymptotic value. Here also the transition functions can be taken as smooth as one wants.

6.4. Example 5.2 + 5.2.

By coupling Example 5.2 with a similar configuration controlled by the other Player, one recovers exactly the family of counterexamples in [12]. Note that in this case both configurations are oscillating, but with a different phase so the ratio does not converge.

6.5. Example 5.4 + 5.4, 5.5 + 5.5 and 5.6 + 5.6.

Two examples of Ziliotto ([13], sections 2.1 2.2 and 4.1) are combinations of either 5.5, 5.6 or 5.7 with a similar configuration. In those cases both configurations are oscillating of order $\frac{1}{2}$ but one is oscillating twice as fast as the other hence the oscillations of v_λ in the combined game.

6.6. Example 4.1 + 5.4.

This gives a MDP with a countable number of states (and only 2 actions) in which v_λ does not converge. Observe that one can compactify the state space in such a way that both the payoff and transition functions are continuous.

7. COMPARISON AND CONCLUSION**7.1. Irreversibility.**

The above analysis shows that oscillations in the inertia rate and reversibility allows for non convergence of the discounted values.

These two properties seem to be also necessary. In fact, Sorin and Vigeral [11] prove the existence of the limit of the discounted values for stochastic games with finite state space, continuous action space and continuous payoffs and transitions for absorbing games see also [5, 6, 9] and recursive games see also [10]. These two classes corresponds to the “irreversible” case where once one leaves a state, it cannot be reached again.

7.2. Remark that any oscillating configuration of Section 5 leads, under optimal play, to an almost immediate exit. Hence, by itself, any such configuration leads to a regular asymptotic behavior. It is only the “resonance” between two configurations that yields asymptotic issues.

7.3. Semi-algebraic.

In the case of stochastic games with finitely many states and full monitoring, in all the examples of the previous section there is a lack of semi-algebraicity, either because transition functions oscillate infinitely often or because a set of actions has infinitely many connected components. While the existence of an asymptotic value with semi-algebraic parameters in the case of either perfect information or finitely many actions on one side holds [2], it is not known in full generality. In particular, an interesting question is to determine whether there exists a configuration with semi-algebraic parameters such that Q_λ is not semi-algebraic.

7.4. Related issues.

The stationarity of the model is crucial here. However it is possible to construct similar examples in which $\lim v_n$ does not exist. The idea grounds on a lemma of Neyman [8] giving sufficient conditions, for the two sequences v_n and v_{λ_n} for $\lambda_n = \frac{1}{n}$, to have the same asymptotic behavior as n tends to infinity. See [12] for specific details in the framework of sections 5.1 and 5.2 and [13] in the framework of section 5.3-5.6.

Acknowledgments

We thank Jerome Renault and Bruno Ziliotto for helpful comments.

REFERENCES

- [1] T. Bewley and E. Kohlberg (1978), On stochastic games with stationary optimal strategies. *Mathematics of Operations Research* **3** 104-125.
- [2] J. Bolte, S. Gaubert and G. Vigeral (2013), Definable zero-sum stochastic games. <http://arxiv.org/abs/1301.1967>.
- [3] G. Grimmett and D. Stirzaker (2001), *Probability and Random Processes*. Oxford university press
- [4] R. Laraki (2001). Variational inequalities, system of functional equations, and incomplete information repeated games. *SIAM J. Control and Optimization*, **40**, 516-524.
- [5] R. Laraki (2010). Explicit formulas for repeated games with absorbing states. *International Journal of Game Theory*, **39**, 53-69.
- [6] J.-F. Mertens, A. Neyman and D. Rosenberg (2009), Absorbing games with compact action spaces. *Mathematics of Operation Research* **34** 257-262.
- [7] J.-F. Mertens, S. Sorin and S. Zamir (1994). *Repeated Games*. CORE DP 9420-22.
- [8] A. Neyman (2003) , Stochastic games and nonexpansive maps. Chapter 26 in A. Neyman and S. Sorin (eds), *Stochastic Games and Applications*, Kluwer Academic Publishers.
- [9] D. Rosenberg and S. Sorin (2001). An operator approach to zero-sum repeated games. *Israel Journal of Mathematics*, **121**, 221-246.
- [10] S. Sorin (2003), The operator approach to zero-sum stochastic games. Chapter 27 in A. Neyman and S. Sorin (eds), *Stochastic Games and Applications*, Kluwer Academic Publishers.
- [11] S. Sorin and G. Vigeral (2013) Existence of the limit value of two person zero-sum discounted repeated games via comparison theorems. *Journal of Opimization Theory and Applications* **157** , 564-576.
- [12] G. Vigeral (2013) A zero-sum stochastic game with compact action sets and no asymptotic value, *Dynamic Games and Applications*, **3**, 172-186.
- [13] B. Ziliotto (2013) Zero-sum repeated games: counterexamples to the existence of the asymptotic value and the conjecture $\max\min = \lim v_n$, hal- 00824039.

Sylvain Sorin, COMBINATOIRE ET OPTIMISATION, IMJ, CNRS UMR 7586, FACULTÉ DE MATHÉMATIQUES, UNIVERSITÉ P. ET M. CURIE - PARIS 6, TOUR 15-16, 1 ÉTAGE, 4 PLACE JUSSIEU, 75005 PARIS, FRANCE

E-mail address: sorin@math.jussieu.fr

<http://www.math.jussieu.fr/~sorin/>

Guillaume Vigeral, UNIVERSITÉ PARIS-DAUPHINE, CEREMADE, PLACE DU MARÉCHAL DE LATTRE DE TASSIGNY. 75775 PARIS CEDEX 16, FRANCE

E-mail address: vigeral@ceremade.dauphine.fr

<http://www.ceremade.dauphine.fr/~vigeral/indexenglish.html>