

# Zero-Sum Stochastic Games with Vanishing Stage Duration and Public Signals

Ivan Novikov

Université Paris-Dauphine, CEREMADE

28/05/2024

## Table of contents

## Zero-sum stochastic games with perfect observation of the state

## Stochastic games with stage duration





## Strategies and total payoff

- Strategies  $\sigma, \tau$  of players consist in choosing at each stage a mixed action;
- The players can take into account the previous actions of players, as well as the current and previous states.
- $\lambda$ -discounted total payoff:  $E_{\sigma, \tau}^{\omega} \left( \lambda \sum_{i=1}^{\infty} (1 - \lambda)^{i-1} g_i \right)$ ;
- Depends on  $\lambda \in (0, 1)$ , initial state  $\omega$ , and strategies of the players;
- Value  $v_{\lambda} : \Omega \rightarrow \mathbb{R}$ :

$$\begin{aligned} v_\lambda(\omega) &= \sup_{\sigma} \inf_{\tau} E_{\sigma, \tau}^{\omega} \left( \lambda \sum_{i=1}^{\infty} (1-\lambda)^{i-1} g_i \right) \\ &= \inf_{\tau} \sup_{\sigma} E_{\sigma, \tau}^{\omega} \left( \lambda \sum_{i=1}^{\infty} (1-\lambda)^{i-1} g_i \right). \end{aligned}$$



## Stochastic games with stage duration

## Kernel

- Kernel  $q : I \times J \times \Omega \rightarrow \mathbb{R}^{|\Omega|}$ .

$$q(i, j, \omega)(\omega') = \begin{cases} P(i, j, \omega)(\omega') & \text{if } \omega \neq \omega'; \\ P(i, j, \omega)(\omega') - 1 & \text{if } \omega = \omega'. \end{cases}$$

- Recall that  $P(i, j, \omega)(\omega')$  is the probability that the next state is  $\omega'$ , if the current state is  $\omega$  and players' actions are  $(i, j)$ ;
- Hence the closer kernel  $q$  is to 0, the more probable it is that the next state coincides with the current one.







## Comparison (1)

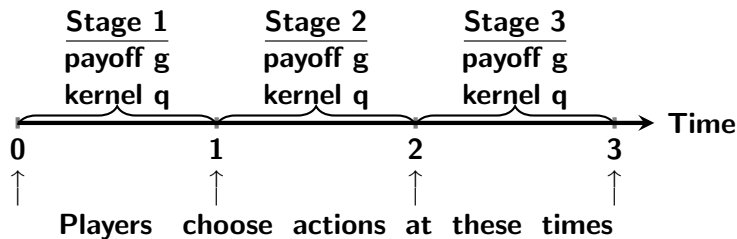


Figure: "Usual" stochastic game: duration of each stage is 1



## Discounted games with stage duration

- For a game with stage duration  $h$ , the total payoff is (depending on the discount factor  $\lambda$ , initial state  $\omega$ , and strategies  $\sigma, \tau$  of players)

$$E_{\sigma,\tau}^\omega\left(\lambda\sum_{k=1}^\infty(1-\lambda h)^{k-1}(g_k)_h\right);$$

- Why such a choice? Easy explanation:
- The total payoff is  $\lambda$ -discounted game with stage duration 1 is  $E_{\sigma, \tau}^{\omega} (\lambda \sum_{k=1}^{\infty} (1 - \lambda)^{k-1} g_k)$ . The total payoff of  $\lambda$ -discounted game with stage duration  $h$  is  $E_{\sigma, \tau}^{\omega} (\sum_{k=1}^{\infty} \lambda h (1 - \lambda h)^{k-1} g_k)$ ;
- So, it may be seen as a game with discount factor  $\lambda h$ . I.e., the discount factor is proportional to  $h$ , just as the payoff  $g$  and the kernel  $q$ .









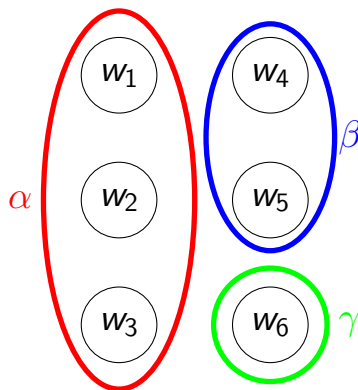


- Now players cannot perfectly observe the current state;
- Players know the initial probability distribution on the states and some information about the current state.





## An example of the partition function $f$



There are 3 public signals, and  $f(w_1) = f(w_2) = f(w_3) = \alpha$ ,  
 $f(w_4) = f(w_5) = \beta$ ,  $f(w_6) = \gamma$ .

- We still can consider games with stage duration  $h$  in this new setting;
- Payoff  $g_h = hg$ ;
- Kernel  $q_h = hq$ ;
- State space  $\Omega$ , signal set  $A$ , partition function  $f$ , and action spaces  $I$  and  $J$  of player 1 and player 2 are independent of  $h$ ;
- The total payoff is still  $E_{\sigma, \tau}^{\omega} \left( \lambda \sum_{k=1}^{\infty} (1 - \lambda h)^{k-1} (g_k)_h \right)$ ;
- $v_{h, \lambda}$  is the value of the game with such a total payoff.





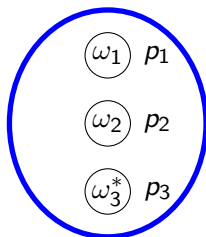




## Second result (1)

### Theorem

*There is a stochastic game  $G$  with public signals in which the uniform limit  $\lim_{\lambda \rightarrow 0} \lim_{h \rightarrow 0} v_{h,\lambda}$  exists, but the pointwise limit  $\lim_{\lambda \rightarrow 0} v_{1,\lambda}$  does not exist.*

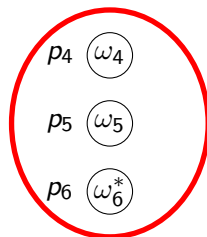


**Signal MINUS**

**Payoff -1**

Player 1's actions:  $T, B, Q$

Player 2's actions:  $L, R$



**Signal PLUS**

**Payoff +1**

Player 1's actions:  $T, M, B$

Player 2's actions:  $L, M, R, Q$

## Second result (2)

The transition matrices for non-absorbing states:

State  $\omega_1$ :

	$L$	$R$
$T$	$\omega_1$	$\omega_2$
$B$	$\omega_2$	$\omega_1$
$Q$	$\omega_5$	$\omega_5$

State  $\omega_2$ :

	$L$	$R$
$T$	$\frac{1}{2}\omega_1 + \frac{1}{2}\omega_2$	$\omega_2$
$B$	$\omega_2$	$\frac{1}{2}\omega_1 + \frac{1}{2}\omega_2$
$Q$	$\omega_3^*$	$\omega_3^*$

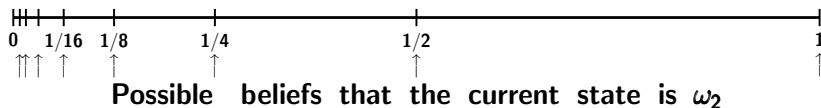
State  $\omega_4$ :

	$L$	$M$	$R$	$Q$
$T$	$\omega_4$	$\omega_5$	$\omega_5$	$\omega_2$
$M$	$\omega_5$	$\omega_4$	$\omega_5$	$\omega_2$
$B$	$\omega_5$	$\omega_5$	$\omega_4$	$\omega_2$

State  $\omega_5$ :

	$L$	$M$	$R$	$Q$
$T$	$\frac{2}{3}\omega_4 + \frac{1}{3}\omega_5$	$\omega_5$	$\omega_5$	$\omega_6^*$
$M$	$\omega_5$	$\frac{2}{3}\omega_4 + \frac{1}{3}\omega_5$	$\omega_5$	$\omega_6^*$
$B$	$\omega_5$	$\omega_5$	$\frac{2}{3}\omega_4 + \frac{1}{3}\omega_5$	$\omega_6^*$

## Informal proof (1)



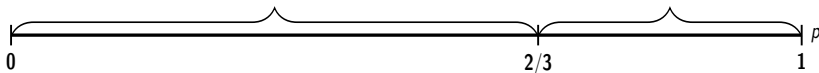
**Figure:** Discrete case (i.e. stage duration is  $h = 1$ ). Possible beliefs of player 1 that the current state is  $\omega_2$  if player 2 plays optimally. As  $\lambda$  becomes smaller, player 1 can wait longer and longer to achieve higher probabilities.

- If the current signal is LEFT, then the smaller is the discount factor  $\lambda$ , the smaller is player 1 can make his belief that the current state is  $\omega_2$ ;
- Analogously, if the current signal is RIGHT, then the smaller is  $\lambda$ , the smaller is player 2 can make his belief that the current state is  $\omega_5$ ;
- Because of that, there is an oscillation when  $\lambda \rightarrow 0$ .

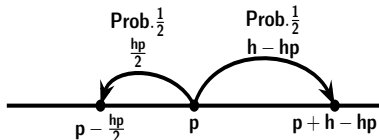
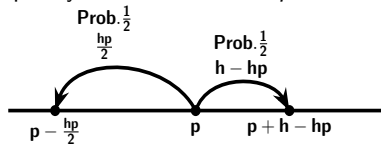
## Informal proof (2)

Player 1 immediately starts playing  $Q$

Player 1 plays  $C$  until it gets sufficiently close to  $p = 2/3$ .



**Figure:** Continuous case (i.e.  $h \approx 0$ ) with small  $\lambda$ . With prob.  $p < 2/3$  that the current state is  $\omega_2$ , player 1 should immediately start playing  $Q$ . Otherwise, his belief  $\tilde{p}$  will start to increase until it becomes  $\tilde{p} = 2/3$ , which is bad for player 1. With prob.  $p \geq 2/3$  that the current state is  $\omega_2$ , player 1 can very quickly decrease his belief  $\tilde{p}$  until it becomes  $\tilde{p} \approx 2/3$ , which is good for him.



(a)  $p > 2/3$  and player 1 plays not  $Q$ .  
 $E(\tilde{p} - p) = \frac{1}{2}(h - hp) + \frac{1}{2} \cdot \frac{-hp}{2} = \frac{h}{4}(2 - 3p) < 0$ , thus if  $\lambda$  is small, then player 1 prefers do not play  $Q$  until  $\tilde{p}$  is close to  $2/3$ .

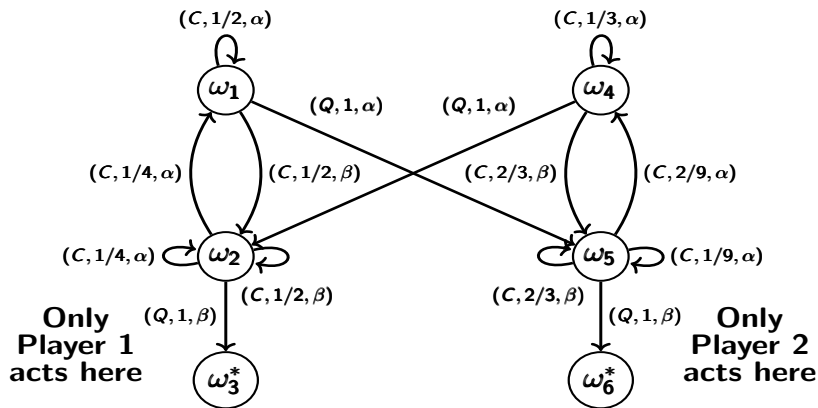
(b)  $p < 2/3$  and player 1 plays not  $Q$ .  
 $E(\tilde{p} - p) = \frac{1}{2}(h - hp) + \frac{1}{2} \cdot \frac{-hp}{2} = \frac{h}{4}(2 - 3p) > 0$ , thus player 1 prefers to play  $Q$  until the state changes.





## The limit $\lim_{\lambda \rightarrow 0} v_{1,\lambda}$ does not exist (2)

Auxiliary game  $\tilde{G}$  with discounted value  $\tilde{v}_\lambda$ .



**Payoff -1**

**Payoff +1**

The arrow from state  $s_1$  to the state  $s_2$  with label  $(X, p, \gamma)$  tells that if player that controls state  $s_1$  chooses action  $X$ , then with probability  $p$  he goes to state  $s_2$  and receives signal  $\gamma$ .



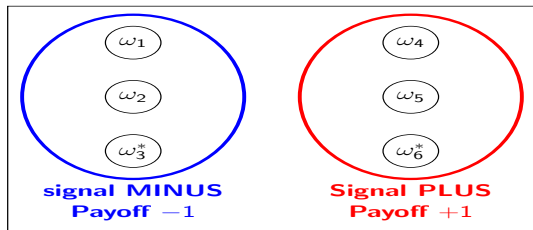
The limit  $\lim_{\lambda \rightarrow 0} v_{1,\lambda}$  does not exist (3)

## Theorem

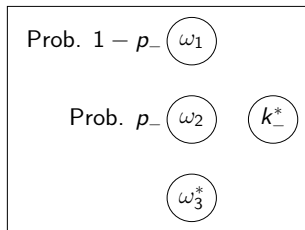
*There is a game in which the uniform limit  $\lim_{\lambda \rightarrow 0} \lim_{h \rightarrow 0} v_{h,\lambda}$  exists, but the pointwise limit  $\lim_{\lambda \rightarrow 0} v_{1,\lambda}$  does not exist.*

- One can prove that  $\tilde{v}_\lambda(p) = v_{1,\lambda}(p)$ . (e.g. by writing the Shapley equation for both games).
- The article “Hidden stochastic games and limit equilibrium payoffs” (Jérôme Renault and Bruno Ziliotto, 2020) proves that  $\lim_{\lambda \rightarrow 0} \tilde{v}_\lambda(p)$  does not exist.

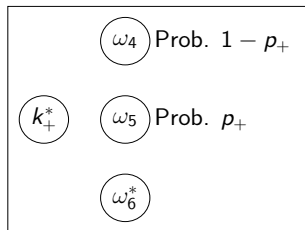
The limit  $\lim_{\lambda \rightarrow 0} \lim_{h \rightarrow 0} v_{h,\lambda}$  exists (1)



↓ Game  $G$  with two public signals ↓



(a) State-blind “half-game”  
 $G^-(k_-)$ , where  $k_- \in [-1, 1]$ .



(b) State-blind “half-game”  
 $G^+(k_+)$ , where  $k_+ \in [-1, 1]$ .

The limit  $\lim_{\lambda \rightarrow 0} \lim_{h \rightarrow 0} v_{h,\lambda}$  exists (2)

- We only need to find the values  $v_{h,\lambda}^-(k, p)$  and  $v_{h,\lambda}^+(k, p)$  of these two “half-games”!
- In this case we can deduce  $\lim_{\lambda \rightarrow 0} \lim_{h \rightarrow 0} v_{h,\lambda}$  for initial states  $\omega_2$  and  $\omega_5$  by solving a system of two equations with variables  $k_-$  and  $k_+$ . We have 
$$\begin{cases} v_{h,\lambda}(\omega_2) = v_{h,\lambda}^-(v_{h,\lambda}(\omega_5), \omega_2) \\ v_{h,\lambda}(\omega_5) = v_{h,\lambda}^+(v_{h,\lambda}(\omega_2), \omega_5) \end{cases}.$$
- Later we can find  $\lim_{\lambda \rightarrow 0} \lim_{h \rightarrow 0} v_{h,\lambda}(p)$  for any initial  $p$  by replacing  $k_-$  or  $k_+$  with values that were just found.

The limit  $\lim_{\lambda \rightarrow 0} \lim_{h \rightarrow 0} v_{h,\lambda}$  exists (3)

Denote by  $v_{h,\lambda}^-(k, p)$  the value of the game  $G^-(k)$  with initial state  $\omega_2$  (prob.  $p$ ) or  $\omega_1$  (prob.  $1 - p$ );

## Lemma

For any  $p \in [0, 1]$  and any  $k \in [-1, 1]$  we have

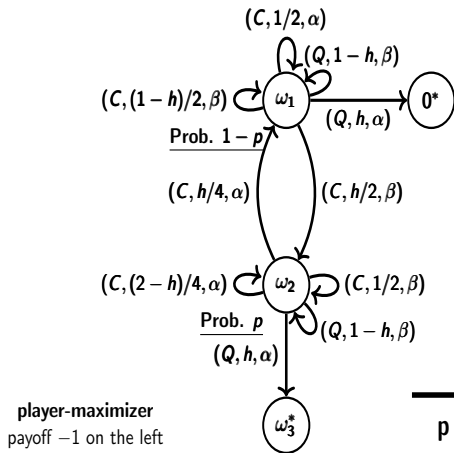
$$v_{h,\lambda}^-(k, p) = (k + 1)v_{h,\lambda}^-(0, p) + k.$$

Proof: Since  $k \geq -1$ , any optimal strategy in the game  $G^-(0)$  is also optimal in the game  $G^-(k)$ . Thus there exists  $\alpha \in [0, 1]$  such that we have  $v_{h,\lambda}^-(k, p) = \alpha k + (-1)(1 - \alpha)$ . By taking  $k = 0$ , we obtain

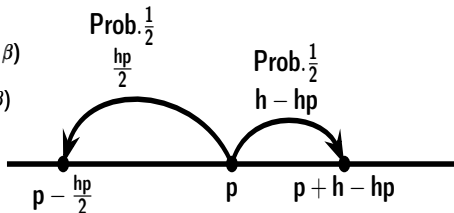
$$v_{h,\lambda}^-(0, p) = -1 + \alpha \iff \alpha = v_{h,\lambda}^-(0, p) + 1.$$



## How to guess $w(p)$ ? (1)



(a) Auxiliary game



(b) Prob.  $p$  that the current state is  $\omega_2$

The arrow from state  $s_1$  to the state  $s_2$  with label  $(X, p, \gamma)$  tells that if player plays  $X$ , then with prob.  $p$  he goes to state  $s_2$  and receives signal  $\gamma$

## How to guess $w(p)$ ? (2)

- This new game is equivalent to  $G^-(0)$ ;
- Denote by  $w(p)$  its value;
- Shapley equation:

$$w(p) = -\lambda h + (1 - \lambda h) \max \left\{ \underbrace{-hp}_{\text{Player 1 plays Q, the signal is } \alpha} + \underbrace{(1 - h)w(p)}_{\text{Player 1 plays Q, the signal is } \beta}; \right. \\ \left. \underbrace{\frac{1}{2}w\left(p - \frac{hp}{2}\right)}_{\text{Player 1 plays C, the signal is } \alpha} + \underbrace{\frac{1}{2}w(p + h - hp)}_{\text{Player 1 plays C, the signal is } \beta} \right\}.$$

- there is  $p^* \in [0, 1]$  such that for  $p \leq p^*$  player prefers to play  $Q$ , and for  $p > p^*$  player prefers to play  $C$ . Thus for  $p \leq p^*$

$$-\lambda h + (1 - \lambda h)(-hp + (1 - h)w(p)) = w(p) \iff$$

$$w(p) = \frac{(h\lambda - 1)p - \lambda}{1 + (1 - h)\lambda} \xrightarrow{h \rightarrow 0} -\frac{p + \lambda}{1 + \lambda}.$$

### How to guess $w(p)$ ? (3)

$p^*$  is an approximate solution of the equation

$$-hp + (1-h)\frac{(h\lambda-1)p-\lambda}{1+(1-h)\lambda} = \frac{(h\lambda-1)\left(p-\frac{hp}{2}\right)-\lambda}{2(1+(1-h)\lambda)} + \frac{(h\lambda-1)(p+h-hp)-\lambda}{2(1+(1-h)\lambda)},$$

from which  $p^* = \frac{4\lambda+2-2\lambda h}{4\lambda+3-7\lambda h} \xrightarrow{h \rightarrow 0} \frac{4\lambda+2}{4\lambda+3}$ .



## How to guess $w(p)$ ? (4)

- For  $p \geq p^*$ ,  $w(p)$  is a solution of the equation (in  $f(p)$ )

$$f(p) = -\lambda h + (1 - \lambda h) \left( \frac{1}{2} f \left( p - \frac{hp}{2} \right) + \frac{1}{2} f(p + h - hp) \right).$$

- if  $w(p)$  is differentiable, then

$$w(p) = -\lambda h + (1 - \lambda h) \left( \frac{1}{2} \left( w(p) - \frac{1}{2} hp w'(p) \right) + \frac{1}{2} (w(p) + (h - hp) w'(p)) \right) + o(h).$$

- Thus we have for small  $h$

$$\begin{cases} \lambda w(p) \approx -\lambda - \frac{1}{4} p w'(p) + \frac{1}{2} (1 - p) w'(p), & \text{if } p \in (p^*, 1); \\ w(p^*) = \frac{-p^* - \lambda}{1 + \lambda}. \end{cases}$$

from which

$$w(p) = -1 + \frac{(4\lambda)^{4\lambda/3}}{(1 + \lambda)(3 + 4\lambda)^{1+(4\lambda/3)}} (3p - 2)^{-4\lambda/3}.$$

## Theorem

*There is a stochastic game  $G$  with public signals in which the uniform limit  $\lim_{\lambda \rightarrow 0} \lim_{h \rightarrow 0} v_{h,\lambda}$  exists, but the pointwise limit  $\lim_{\lambda \rightarrow 0} v_{1,\lambda}$  does not exist.*

Open question: For the considered above game  $G$ , can we say that

1. For any fixed  $h \in (0, 1]$ , the limit  $\lim_{\lambda \rightarrow 0} v_{h,\lambda}$  does not exist?
2. We have  $\left| \limsup_{\lambda \rightarrow 0} v_{h,\lambda}(p) - \liminf_{\lambda \rightarrow 0} v_{h,\lambda}(p) \right| \rightarrow 0$  as  $h \rightarrow 0$ , uniformly in  $p$ ?



# Thank you!