

## Repeated and Stochastic Games

Branislav Bošanský

Artificial Intelligence Center,  
Department of Computer Science,  
Faculty of Electrical Engineering,  
Czech Technical University in Prague

*branislav.bosansky@agents.fel.cvut.cz*

April 1, 2019

# Repeated Games

# Repeated Games

Repeated Games are the simplest type of a dynamic game that evolves over time.

# Repeated Games

Repeated Games are the simplest type of a dynamic game that evolves over time.

As such we can treat them as an extensive-form game (the finitely repeated case), or a stochastic game (the infinitely repeated case). However, such representations are very inefficient.

# Repeated Games

Repeated Games are the simplest type of a dynamic game that evolves over time.

As such we can treat them as an extensive-form game (the finitely repeated case), or a stochastic game (the infinitely repeated case). However, such representations are very inefficient.

Repeated games can thus be seen as an example of a compact representation.

# Repeated Games

Repeated Games are the simplest type of a dynamic game that evolves over time.

As such we can treat them as an extensive-form game (the finitely repeated case), or a stochastic game (the infinitely repeated case). However, such representations are very inefficient.

Repeated games can thus be seen as an example of a compact representation.

	<i>C</i>	<i>D</i>
<i>C</i>	(1, 1)	(-1, 2)
<i>D</i>	(2, -1)	(0, 0)

# Repeated Games

Repeated Games are the simplest type of a dynamic game that evolves over time.

As such we can treat them as an extensive-form game (the finitely repeated case), or a stochastic game (the infinitely repeated case). However, such representations are very inefficient.

Repeated games can thus be seen as an example of a compact representation.

	<i>C</i>	<i>D</i>
<i>C</i>	(1, 1)	(-1, 2)
<i>D</i>	(2, -1)	(0, 0)

Natural question: Is a NE of a single game the same as in the (in)finitely repeated game?

# Repeated Games



# Repeated Games

## Definition

Let  $G' = (\mathcal{N}, \mathcal{A}, u)$  be a normal-form game. An **infinitely repeated game** with discounted payoff is an extensive-form game with simultaneous moves  $G^\infty = (\mathcal{N}, \mathcal{H}, \mathcal{A}, g, \delta)$ , where

# Repeated Games

## Definition

Let  $G' = (\mathcal{N}, \mathcal{A}, u)$  be a normal-form game. An **infinitely repeated game** with discounted payoff is an extensive-form game with simultaneous moves  $G^\infty = (\mathcal{N}, \mathcal{H}, \mathcal{A}, g, \delta)$ , where

- $\mathcal{H} = \{\emptyset\} \cup \bigcup_{t=1}^{\infty} A^t \cup A^\infty$

# Repeated Games

## Definition

Let  $G' = (\mathcal{N}, \mathcal{A}, u)$  be a normal-form game. An **infinitely repeated game** with discounted payoff is an extensive-form game with simultaneous moves  $G^\infty = (\mathcal{N}, \mathcal{H}, \mathcal{A}, g, \delta)$ , where

- $\mathcal{H} = \{\emptyset\} \cup \bigcup_{t=1}^{\infty} A^t \cup A^\infty$
- $\mathcal{S}_i : \mathcal{H} \rightarrow \mathcal{A}_i$

# Repeated Games

## Definition

Let  $G' = (\mathcal{N}, \mathcal{A}, u)$  be a normal-form game. An **infinitely repeated game** with discounted payoff is an extensive-form game with simultaneous moves  $G^\infty = (\mathcal{N}, \mathcal{H}, \mathcal{A}, g, \delta)$ , where

- $\mathcal{H} = \{\emptyset\} \cup \bigcup_{t=1}^{\infty} A^t \cup A^\infty$
- $S_i : \mathcal{H} \rightarrow \mathcal{A}_i$
- $g_i(s_i, s_{-i}) = (1 - \delta) \sum_{t=1}^{\infty} \delta^t \mathbb{E}_{a_i \sim s_i, a_{-i} \sim s_{-i}} (u_i(a_i, a_{-i}))$

# Repeated Games

## Definition

Let  $G' = (\mathcal{N}, \mathcal{A}, u)$  be a normal-form game. An **infinitely repeated game** with discounted payoff is an extensive-form game with simultaneous moves  $G^\infty = (\mathcal{N}, \mathcal{H}, \mathcal{A}, g, \delta)$ , where

- $\mathcal{H} = \{\emptyset\} \cup \bigcup_{t=1}^{\infty} A^t \cup A^\infty$
- $\mathcal{S}_i : \mathcal{H} \rightarrow \mathcal{A}_i$
- $g_i(s_i, s_{-i}) = (1 - \delta) \sum_{t=1}^{\infty} \delta^t \mathbb{E}_{a_i \sim s_i, a_{-i} \sim s_{-i}} (u_i(a_i, a_{-i}))$
- $\delta \in (0, 1)$  is the discount factor

# Repeated Games

# Repeated Games

We can define alternative utility functions in repeated games based on payoff vectors  $v_i^t$  for each:

# Repeated Games

We can define alternative utility functions in repeated games based on payoff vectors  $v_i^t$  for each:

- overtaking payoff:  $\lim_{T \rightarrow \infty} \sum_{t=1}^T v_i^t$



# Repeated Games

We can define alternative utility functions in repeated games based on payoff vectors  $v_i^t$  for each:

- overtaking payoff:  $\lim_{T \rightarrow \infty} \sum_{t=1}^T v_i^t$
- average payoff (or limit mean payoff):  $\lim_{T \rightarrow \infty} \sum_{t=1}^T v_i^t / T$

# Repeated Games

We can define alternative utility functions in repeated games based on payoff vectors  $v_i^t$  for each:

- overtaking payoff:  $\lim_{T \rightarrow \infty} \sum_{t=1}^T v_i^t$
- average payoff (or limit mean payoff):  $\lim_{T \rightarrow \infty} \sum_{t=1}^T v_i^t / T$

## Definition

Player  $i$ 's min-max payoff is

$$\underline{v}_i = \min_{s_{-i}} \max_{s_i} g_i(s_i, s_{-i})$$

# Repeated Games

We can define alternative utility functions in repeated games based on payoff vectors  $v_i^t$  for each:

- overtaking payoff:  $\lim_{T \rightarrow \infty} \sum_{t=1}^T v_i^t$
- average payoff (or limit mean payoff):  $\lim_{T \rightarrow \infty} \sum_{t=1}^T v_i^t / T$

## Definition

Player  $i$ 's min-max payoff is

$$\underline{v}_i = \min_{s_{-i}} \max_{s_i} g_i(s_i, s_{-i})$$

A strategy  $s$  is *individually rational* if  $g_i(s) \geq \underline{v}_i$

# Repeated Games

# Repeated Games

## Theorem (Nash Folk Theorem)

*If  $v_i$  is a feasible and an individually rational payoff, then there exists a discount factor  $\underline{\delta} < 1$  such that for all  $\delta > \underline{\delta}$ , there is a Nash equilibrium of  $G$  with payoff  $v_i$ .*

# Repeated Games

## Theorem (Nash Folk Theorem)

*If  $v_i$  is a feasible and an individually rational payoff, then there exists a discount factor  $\underline{\delta} < 1$  such that for all  $\delta > \underline{\delta}$ , there is a Nash equilibrium of  $G$  with payoff  $v_i$ .*

## Proof.

If  $v_i$  is feasible then there exist a strategy  $s$  such that  $g_i(s) = v_i$  and let  $m_{-i}$  be the minmax strategy of other players to reach value  $\underline{v}_i$  for player  $i$ . Let consider the following strategy:

# Repeated Games

## Theorem (Nash Folk Theorem)

*If  $v_i$  is a feasible and an individually rational payoff, then there exists a discount factor  $\underline{\delta} < 1$  such that for all  $\delta > \underline{\delta}$ , there is a Nash equilibrium of  $G$  with payoff  $v_i$ .*

## Proof.

If  $v_i$  is feasible then there exist a strategy  $s$  such that  $g_i(s) = v_i$  and let  $m_{-i}$  be the minmax strategy of other players to reach value  $\underline{v}_i$  for player  $i$ . Let consider the following strategy:

- 1 play according to  $s_i$  as long as no one deviates

# Repeated Games

## Theorem (Nash Folk Theorem)

*If  $v_i$  is a feasible and an individually rational payoff, then there exists a discount factor  $\underline{\delta} < 1$  such that for all  $\delta > \underline{\delta}$ , there is a Nash equilibrium of  $G$  with payoff  $v_i$ .*

## Proof.

If  $v_i$  is feasible then there exist a strategy  $s$  such that  $g_i(s) = v_i$  and let  $m_{-i}$  be the minmax strategy of other players to reach value  $\underline{v}_i$  for player  $i$ . Let consider the following strategy:

- 1 play according to  $s_i$  as long as no one deviates
- 2 let  $\bar{v}_i$  be the maximum value player  $i$  can get by a deviation in step  $t$



# Repeated Games

## Theorem (Nash Folk Theorem)

*If  $v_i$  is a feasible and an individually rational payoff, then there exists a discount factor  $\underline{\delta} < 1$  such that for all  $\delta > \underline{\delta}$ , there is a Nash equilibrium of  $G$  with payoff  $v_i$ .*

## Proof.

If  $v_i$  is feasible then there exist a strategy  $s$  such that  $g_i(s) = v_i$  and let  $m_{-i}$  be the minmax strategy of other players to reach value  $\underline{v}_i$  for player  $i$ . Let consider the following strategy:

- 1 play according to  $s_i$  as long as no one deviates
- 2 let  $\bar{v}_i$  be the maximum value player  $i$  can get by a deviation in step  $t$

$$\begin{aligned}(1 - \delta)[v_i + \delta v_i + \dots + \delta^t \bar{v}_i + \delta^{t+1} \underline{v}_i + \dots] &\leq \\ &\leq (1 - \delta)[v_i + \delta v_i + \dots + \delta^t v_i + \delta^{t+1} v_i + \dots]\end{aligned}$$

# Repeated Games

(Proof cont.)

By setting  $\underline{\delta}$  sufficiently large approaching 1 the above inequality holds. □

# Repeated Games

(Proof cont.)

By setting  $\underline{\delta}$  sufficiently large approaching 1 the above inequality holds. □

The Nash folk theorem says that essentially anything goes as a Nash equilibrium payoff in a discounted repeated game.

# Repeated Games

(Proof cont.)

By setting  $\underline{\delta}$  sufficiently large approaching 1 the above inequality holds. □

The Nash folk theorem says that essentially anything goes as a Nash equilibrium payoff in a discounted repeated game.

The players threat by playing *grim trigger* strategies.

# Stochastic Games

# Stochastic Games

Let's generalize the repeated games.

# Stochastic Games

Let's generalize the repeated games. We do not have to play the same normal-form game repeatedly.

# Stochastic Games

Let's generalize the repeated games. We do not have to play the same normal-form game repeatedly. We can play different normal-form games (possibly for infinitely long time).



# Stochastic Games

Let's generalize the repeated games. We do not have to play the same normal-form game repeatedly. We can play different normal-form games (possibly for infinitely long time).

## Definition (Stochastic game)

A *stochastic game* is a tuple  $(Q, \mathcal{N}, \mathcal{A}, \mathcal{P}, \mathcal{R})$ , where:

# Stochastic Games

Let's generalize the repeated games. We do not have to play the same normal-form game repeatedly. We can play different normal-form games (possibly for infinitely long time).

## Definition (Stochastic game)

A *stochastic game* is a tuple  $(Q, \mathcal{N}, \mathcal{A}, \mathcal{P}, \mathcal{R})$ , where:

$Q$  is a finite set of games

# Stochastic Games

Let's generalize the repeated games. We do not have to play the same normal-form game repeatedly. We can play different normal-form games (possibly for infinitely long time).

## Definition (Stochastic game)

A *stochastic game* is a tuple  $(Q, \mathcal{N}, \mathcal{A}, \mathcal{P}, \mathcal{R})$ , where:

$Q$  is a finite set of games

$\mathcal{N}$  is a finite set of players

# Stochastic Games

Let's generalize the repeated games. We do not have to play the same normal-form game repeatedly. We can play different normal-form games (possibly for infinitely long time).

## Definition (Stochastic game)

A *stochastic game* is a tuple  $(Q, \mathcal{N}, \mathcal{A}, \mathcal{P}, \mathcal{R})$ , where:

$Q$  is a finite set of games

$\mathcal{N}$  is a finite set of players

$\mathcal{A}$  is a finite set of actions,  $\mathcal{A}_i$  are actions available to player  $i$

# Stochastic Games

Let's generalize the repeated games. We do not have to play the same normal-form game repeatedly. We can play different normal-form games (possibly for infinitely long time).

## Definition (Stochastic game)

A *stochastic game* is a tuple  $(Q, \mathcal{N}, \mathcal{A}, \mathcal{P}, \mathcal{R})$ , where:

$Q$  is a finite set of games

$\mathcal{N}$  is a finite set of players

$\mathcal{A}$  is a finite set of actions,  $\mathcal{A}_i$  are actions available to player  $i$

$\mathcal{P}$  is a transition function  $\mathcal{P} : Q \times \mathcal{A} \times Q \rightarrow [0, 1]$ , where  $\mathcal{P}(q, a, q')$  is a probability of reaching game  $q'$  after a joint action  $a$  is played in game  $q$

# Stochastic Games

Let's generalize the repeated games. We do not have to play the same normal-form game repeatedly. We can play different normal-form games (possibly for infinitely long time).

## Definition (Stochastic game)

A *stochastic game* is a tuple  $(Q, \mathcal{N}, \mathcal{A}, \mathcal{P}, \mathcal{R})$ , where:

$Q$  is a finite set of games

$\mathcal{N}$  is a finite set of players

$\mathcal{A}$  is a finite set of actions,  $\mathcal{A}_i$  are actions available to player  $i$

$\mathcal{P}$  is a transition function  $\mathcal{P} : Q \times \mathcal{A} \times Q \rightarrow [0, 1]$ , where  $\mathcal{P}(q, a, q')$  is a probability of reaching game  $q'$  after a joint action  $a$  is played in game  $q$

$\mathcal{R}$  is a set of reward functions  $r_i : Q \times \mathcal{A} \rightarrow \mathbb{R}$

# Stochastic Games

# Stochastic Games

Similarly to repeated games we can have several different rewards (or objectives):



# Stochastic Games

Similarly to repeated games we can have several different rewards (or objectives):

- discounted

# Stochastic Games

Similarly to repeated games we can have several different rewards (or objectives):

- discounted
- average

# Stochastic Games

Similarly to repeated games we can have several different rewards (or objectives):

- discounted
- average
- reachability/safety

# Stochastic Games

Similarly to repeated games we can have several different rewards (or objectives):

- discounted
- average
- reachability/safety

In reachability objectives a player wants to visit certain games infinitely often.

# Stochastic Games

Similarly to repeated games we can have several different rewards (or objectives):

- discounted
- average
- reachability/safety

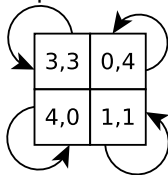
In reachability objectives a player wants to visit certain games infinitely often.

Related to reaching some target state (for example attacking a target) in a game without a pre-determined horizon.

# Stochastic Games - Examples

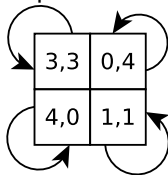
# Stochastic Games - Examples

Repeated prisoners dilemma:

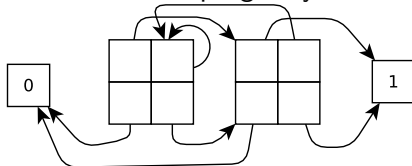


# Stochastic Games - Examples

Repeated prisoners dilemma:



Dante's purgatory:





# Equilibria in Stochastic Games

# Equilibria in Stochastic Games

## Definition (History)

Let  $h_t = (q_0, a_0, q_1, a_1, \dots, a_{t-1}, q_t)$  denote a history of  $t$  stages of a stochastic game, and let  $H_t$  be the set of all possible histories of this length.

# Equilibria in Stochastic Games

## Definition (History)

Let  $h_t = (q_0, a_0, q_1, a_1, \dots, a_{t-1}, q_t)$  denote a history of  $t$  stages of a stochastic game, and let  $H_t$  be the set of all possible histories of this length.

## Definition (Behavioral strategy)

A behavioral strategy  $s_i(h_t, a_{i_j})$  returns the probability of playing action  $a_{i_j}$  for history  $h_t$ .

# Equilibria in Stochastic Games

## Definition (History)

Let  $h_t = (q_0, a_0, q_1, a_1, \dots, a_{t-1}, q_t)$  denote a history of  $t$  stages of a stochastic game, and let  $H_t$  be the set of all possible histories of this length.

## Definition (Behavioral strategy)

A behavioral strategy  $s_i(h_t, a_{i_j})$  returns the probability of playing action  $a_{i_j}$  for history  $h_t$ .

## Definition (Markov strategy)

A Markov strategy  $s_i$  is a behavioral strategy in which  $s_i(h_t, a_{i_j}) = s_i(h'_t, a_{i_j})$  if  $q_t = q'_t$ , where  $q_t$  and  $q'_t$  are the final games of  $h_t$  and  $h'_t$ , respectively.

# Equilibria in Stochastic Games

# Equilibria in Stochastic Games

## Definition

A strategy profile is called a *Markov perfect equilibrium* if it consists of only Markov strategies, and is a Nash equilibrium.

# Equilibria in Stochastic Games

## Definition

A strategy profile is called a *Markov perfect equilibrium* if it consists of only Markov strategies, and is a Nash equilibrium.

## Theorem

*Every  $n$ -player, general-sum, discounted-reward stochastic game has a Markov perfect equilibrium.*

# Equilibria in Stochastic Games

## Definition

A strategy profile is called a *Markov perfect equilibrium* if it consists of only Markov strategies, and is a Nash equilibrium.

## Theorem

*Every  $n$ -player, general-sum, discounted-reward stochastic game has a Markov perfect equilibrium.*

## Theorem

*Problem of computing an optimal strategy in simple (turn-taking) stochastic games, where pure stationary strategies are known to be optimal, is in PLS.*

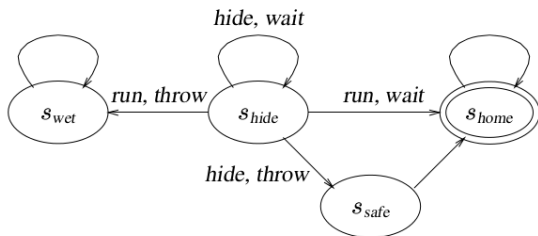
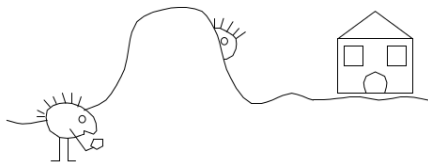


# Equilibria in Stochastic Games

For other rewards, Markov perfect equilibrium does not have to exist.

# Equilibria in Stochastic Games

For other rewards, Markov perfect equilibrium does not have to exist.



# Approximating Optimal Strategies in Stochastic Games

Standard algorithms from Markov Decision Processes, value and strategy iteration, translate to stochastic games.

---

<sup>1</sup>Pseudocode from [3].

# Approximating Optimal Strategies in Stochastic Games

Standard algorithms from Markov Decision Processes, value and strategy iteration, translate to stochastic games.

---

## Algorithm 1. Value Iteration

---

```
1:  $t := 0$ 
2:  $\tilde{v}^0 := (0, \dots, 0, 1)$  // the vector  $\tilde{v}^0$  is indexed  $0, 1, \dots, N, N + 1$ 
3: while true do
4:    $t := t + 1$ 
5:    $\tilde{v}_0^t := 0$ 
6:    $\tilde{v}_{N+1}^t := 1$ 
7:   for  $i \in \{1, 2, \dots, N\}$  do
8:      $\tilde{v}_i^t := \text{val}(A_i(\tilde{v}^{t-1}))$ 
```

---

1

---

<sup>1</sup>Pseudocode from [3].

# Approximating Optimal Strategies in Stochastic Games

---

## Algorithm 2. Strategy Iteration

---

```
1:  $t := 1$ 
2:  $x^1 :=$  the strategy for Player I playing uniformly at each position
3: while true do
4:    $y^t :=$  an optimal best reply by Player II to  $x^t$ 
5:   for  $i \in \{0, 1, 2, \dots, N, N + 1\}$  do
6:      $v_i^t := \mu_i(x^t, y^t)$ 
7:    $t := t + 1$ 
8:   for  $i \in \{1, 2, \dots, N\}$  do
9:     if  $\text{val}(A_i(v^{t-1})) > v_i^{t-1}$  then
10:       $x_i^t := \text{maximin}(A_i(v^{t-1}))$ 
11:     else
12:       $x_i^t := x_i^{t-1}$ 
```

# Stochastic Games with Imperfect Information

Extending stochastic games to imperfect information (known as *partial observability*, hence termed Partially Observable Stochastic Games (POSGs)) is lot more complicated compared to finite EFGs.

# Stochastic Games with Imperfect Information

Extending stochastic games to imperfect information (known as *partial observability*, hence termed Partially Observable Stochastic Games (POSGs)) is lot more complicated compared to finite EFGs.

The problem lies with *Nested beliefs*. Consider a two-player game where each player has some private state unobserved by the opponent:

# Stochastic Games with Imperfect Information

Extending stochastic games to imperfect information (known as *partial observability*, hence termed Partially Observable Stochastic Games (POSGs)) is lot more complicated compared to finite EFGs.

The problem lies with *Nested beliefs*. Consider a two-player game where each player has some private state unobserved by the opponent:

- A player  $i$  has uncertainty about the exact state of the opponent  $-i$  – there is a belief (a probability distribution) over possible states.



# Stochastic Games with Imperfect Information

Extending stochastic games to imperfect information (known as *partial observability*, hence termed Partially Observable Stochastic Games (POSGs)) is lot more complicated compared to finite EFGs.

The problem lies with *Nested beliefs*. Consider a two-player game where each player has some private state unobserved by the opponent:

- A player  $i$  has uncertainty about the exact state of the opponent  $-i$  – there is a belief (a probability distribution) over possible states.
- The optimal strategy of player  $i$  depends on the strategy of the opponent  $-i$  that depends on the belief over possible private states of player  $i$ .

# Stochastic Games with Imperfect Information

Extending stochastic games to imperfect information (known as *partial observability*, hence termed Partially Observable Stochastic Games (POSGs)) is lot more complicated compared to finite EFGs.

The problem lies with *Nested beliefs*. Consider a two-player game where each player has some private state unobserved by the opponent:

- A player  $i$  has uncertainty about the exact state of the opponent  $-i$  – there is a belief (a probability distribution) over possible states.
- The optimal strategy of player  $i$  depends on the strategy of the opponent  $-i$  that depends on the belief over possible private states of player  $i$ .
- Each player needs to consider beliefs, belief of beliefs, ... etc.

# Partially Observable Stochastic Games (POSGs)

Solving general POSGs is not tractable (even solving related single-player decision problems is often undecidable [2]).

# Partially Observable Stochastic Games (POSGs)

Solving general POSGs is not tractable (even solving related single-player decision problems is often undecidable [2]).

We can restrict to subclasses of games with limited partial observability:

# Partially Observable Stochastic Games (POSGs)

Solving general POSGs is not tractable (even solving related single-player decision problems is often undecidable [2]).

We can restrict to subclasses of games with limited partial observability:

- One-Sided Partially Observable Stochastic Games [4]

# Partially Observable Stochastic Games (POSGs)

Solving general POSGs is not tractable (even solving related single-player decision problems is often undecidable [2]).

We can restrict to subclasses of games with limited partial observability:

- One-Sided Partially Observable Stochastic Games [4]
- Partially Observable Stochastic Games with Public Observations [5]

# Partially Observable Stochastic Games (POSGs)

Theoretical results:

# Partially Observable Stochastic Games (POSGs)

Theoretical results:

- Value function (function that assigns a belief point value of a (sub)-game) is convex (or convex-concave).



# Partially Observable Stochastic Games (POSGs)

Theoretical results:

- Value function (function that assigns a belief point value of a (sub)-game) is convex (or convex-concave).
- We can define dynamic-programming operator – a generalization of Bellman update.

$$[Hv](b) = \min_{\pi_2} \max_{\pi_1} \left( R_{\pi_1, \pi_2}^{\text{imm}} + \gamma \cdot R_{\pi_1, \pi_2}^{\text{subs}}(v) \right)$$

# Partially Observable Stochastic Games (POSGs)

Theoretical results:

- Value function (function that assigns a belief point value of a (sub)-game) is convex (or convex-concave).
- We can define dynamic-programming operator – a generalization of Bellman update.

$$[Hv](b) = \min_{\pi_2} \max_{\pi_1} \left( R_{\pi_1, \pi_2}^{\text{imm}} + \gamma \cdot R_{\pi_1, \pi_2}^{\text{subs}}(v) \right)$$

$$R_{\pi_1, \pi_2}^{\text{imm}} = \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}_1} \sum_{a' \in \mathcal{A}_2} b(s) \cdot \pi_1(a) \cdot \pi_2(s, a') \cdot \mathcal{R}(s, a, a')$$

$$R_{\pi_1, \pi_2}^{\text{subs}}(v) = \sum_{a \in \mathcal{A}_1} \sum_{o \in \mathcal{O}} \pi_1(a) \cdot \Pr[o|a, \pi_2] \cdot v(b_{\pi_2}^{a,o})$$

# Partially Observable Stochastic Games (POSGs)

We can generalize value-iteration algorithms for POMDPs to POSGs.

## Heuristic Search Value Iteration (HSVI):

**Data:** Game  $\langle \mathcal{S}, \mathcal{A}_1, \mathcal{A}_2, \mathcal{O}, \mathcal{T}, \mathcal{R} \rangle$ , initial belief  $b^0$ ,  
discount factor  $\gamma$ , desired precision  $\epsilon > 0$ ,  
neighborhood parameter  $R$

**Result:** Approximate value function  $\hat{v}$

```
1 Initialize  $\hat{v}$ 
2 while  $\text{gap}(\hat{v}(b^0)) > \epsilon$  do
3   | Explore( $b^0, \epsilon, R, 0$ )
4 return  $\hat{v}$ 
5 procedure Explore( $b, \epsilon, R, t$ )
6   |  $\pi_2 \leftarrow$  optimal strategy of player 2 in  $[H\underline{v}](b)$ 
7   |  $(a, o) \leftarrow$  select according to forward exploration
   | heuristic
8   | if  $\text{excess}(\hat{v}(b_{\pi_2}^{a,o}), t + 1) > 0$  then
9     | Explore( $b_{\pi_2}^{a,o}, \epsilon, R, t + 1$ )
10  |  $\Gamma \leftarrow \Gamma \cup \{L\underline{\Gamma}(b)\}$ 
11  |  $\Upsilon \leftarrow \Upsilon \cup \{U\underline{\Upsilon}(b)\}$  and make  $\bar{v}$   $(U - L)$ -Lipschitz
```

**Algorithm 1:** HSVI algorithm for one-sided POSGs

# References I

(besides the books)

- [1] M. Osborne and A. Rubinstein, *A course in game theory*. MIT press, 1994.
- [2] O. Madani, S. Hanks, and A. Condon, “On the undecidability of probabilistic planning and infinite-horizon partially observable Markov decision problems,” in *AAAI/IAAI*, pp. 541–548, 1999.
- [3] K. A. Hansen, R. Ibsen-Jensen, and P. B. Miltersen, “The Complexity of Solving Reachability Games Using Value and Strategy Iteration,” in *International Computer Science Symposium in Russia*, pp. 77–90, 2011.
- [4] Horák, K., Božanský, B., and Pěchouček, M. (2017). Heuristic Search Value Iteration for One-Sided Partially Observable Stochastic Games. In *In Proceedings of AAAI Conference on Artificial Intelligence*, pages 558–564.

- [5] Horák, K., Bošanský, B. (2019).  
Solving Partially Observable Stochastic Games with Public Observations.  
In *In Proceedings of AAAI Conference on Artificial Intelligence*.