

# POMDPs and Blind MDPs: (Dis)continuity of Values and Strategies



K. Chatterjee<sup>1</sup>

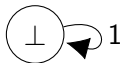


R. Saona<sup>1</sup>

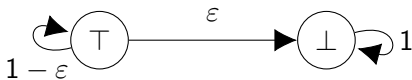
<sup>1</sup>Institute of Science and Technology Austria (ISTA)



# Continuity in Stochastic dynamics



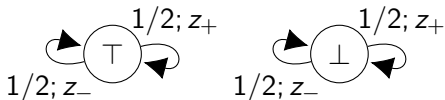
(Deterministic) (dynamic)



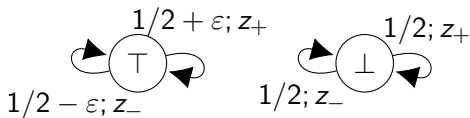
Similar? stochastic dynamic

Stochastic dynamics (MCs) must consider structure when analyzing continuity.

# Continuity in Partially Observable Stochastic dynamics



(Static) partially observable (stochastic) dynamic



Similar? partially observable stochastic dynamic

Belief dynamics are fragile to structurally preserving changes.

# Continuity concepts

- Value-continuity  
Value of similar POMDPs is close
- Weak strategy-continuity  
**Some** approximately-optimal strategy is still approximately-optimal in similar POMDPs
- Strong strategy-continuity  
**All** approximately-optimal strategies are approximately-optimal in similar POMDPs

# Results

Model	Continuity		
	Value	Weak strategy	Strong strategy
Fully-observable MDPs	Yes	Yes	<b>No</b>
POMDPs	<b>No</b>	<b>No</b>	<b>No</b>
Blind MDPs	<b>Yes</b>	<b>Yes</b>	<b>Yes</b>

Theorem: Deciding whether a POMDP is continuous is **algorithmically impossible**.

## Remarks

- Blind MDPs are strictly more well-behaved than POMDPs
- Blind MDPs are strictly more well-behaved than MDPs

A Partially-Observable Markov Decision Process (POMDP) is a tuple  $\Gamma = (\text{States}, \text{Actions}, \text{Signals}, p_1, \delta)$  where

- States is a finite set of **states**;
- Actions is a finite set of **actions**;
- Signals is a finite set of **signals**;
- $p_1 \in \Delta(\text{States})$  is an **initial distribution**;
- $\delta: \text{States} \times \text{Actions} \rightarrow \Delta(\text{States} \times \text{Signals})$  is a probabilistic transition function.

Special cases:

$$\begin{aligned} |\text{Signals}| = 1 &\Rightarrow \text{blind MDP} \\ \text{signal} = \text{state} &\Rightarrow \text{(fully-observable) MDP} \end{aligned}$$

- **strategy**  $\sigma: \bigcup_{n \geq 0} (\text{Actions} \times \text{Signals})^n \rightarrow \Delta(\text{Actions})$
- **play**  $\omega = (s_n, a_n, z_{n+1})_{n \geq 1} \subseteq \text{States} \times \text{Actions} \times \text{Signals}$
- probability  $\mathbb{P}_{\rho_1}^\sigma[\Gamma]$  and expectation  $\mathbb{E}_{\rho_1}^\sigma[\Gamma]$
- **belief**

$$\mathbb{P}_{\rho_1}^\sigma(S_m = \cdot \mid \forall i \in [m-1] \quad A_i = a_i, Z_{i+1} = z_{i+1}),$$

- **reward** reward: States  $\times$  Actions  $\rightarrow \mathbb{R}$
- **objective** payoff( $\omega$ ) is one of

$$\liminf_{m \rightarrow \infty} \frac{1}{m} \sum_{i=1}^m \text{reward}(s_i, a_i) \qquad \limsup_{m \rightarrow \infty} \frac{1}{m} \sum_{i=1}^m \text{reward}(s_i, a_i)$$

$$\liminf_{m \rightarrow \infty} \text{reward}(s_m, a_m)$$

$$\limsup_{m \rightarrow \infty} \text{reward}(s_m, a_m)$$



- **value**

$$\text{val}(\Gamma) := \sup_{\sigma} \mathbb{E}_{p_1}^{\sigma}(\text{payoff}(\omega))$$

- **$\varepsilon$ -optimal strategy**  $\mathbb{E}_{p_1}^{\sigma}(\text{payoff}(\omega)) \geq \text{val}(\Gamma) - \varepsilon$

Special concepts

- **structural equivalence**

$$\text{supp}(\delta(s, a)) = \text{supp}(\tilde{\delta}(s, a))$$

- **$\xi$ -similar POMDPs** For all  $s, a, s', z$ ,

$$\frac{1}{1 + \xi} \delta(s, a)(s', z) \leq \tilde{\delta}(s, a)(s', z) \leq (1 + \xi) \delta(s, a)(s', z)$$

# Results

Model	Continuity		
	Value	Weak strategy	Strong strategy
Fully-observable MDPs	Yes	Yes	<b>No</b>
POMDPs	<b>No</b>	<b>No</b>	<b>No</b>
Blind MDPs	<b>Yes</b>	<b>Yes</b>	<b>Yes</b>

Theorem: Deciding whether a POMDP is value-, weakly strategy-, or strongly strategy-continuous is **algorithmically impossible**.

Blind MDPs:  
no signals  
guarantees continuity

Theorem (Stability of invariant distribution, O’Cinneide 1993)

*Consider an irreducible stochastic matrix  $\Delta$ .*

*Computing the stable distribution*

$$p^\top = p^\top \Delta$$

*is a stable operation.*

The proof is by induction on the dimension of  $\Delta$ , possible thanks to a characterization of the limit

Theorem (Stability of discounted occupation times, Solan 2003)

*Consider a Markov Chain with a starting state.*

*The ratio of  $\lambda$ -discounted occupation times at a state  $s$  is a rational function of the transition probabilities, i.e., for  $\lambda > 0$*

$$\delta \mapsto \frac{\text{time}_\lambda(s, \delta)}{\text{time}_\lambda(s, \tilde{\delta})} = \frac{\text{poly}_s(\delta, \tilde{\delta})}{\text{poly}_s(\delta, \tilde{\delta})}.$$

From this result, we conclude value- and weak strategy-continuity for (fully-observable) MDPs (and zero-sum stochastic games).

# Blind MDPs: Belief dynamic

The belief update in blind MDPs is directly given by the transition.  
For each action  $a$ , define the matrix

$$(M_a)_{s,s'} := \delta(s, a)(s').$$

After playing actions  $a, b, a, \dots$ , the beliefs are

$$p_1^\top \quad p_1^\top M_a \quad p_1^\top M_a M_b \quad p_1^\top M_a M_b M_a \quad \dots$$

For similar matrices  $\tilde{M}_a$ , the beliefs in the corresponding similar blind MDP are

$$p_1^\top \quad p_1^\top \tilde{M}_a \quad p_1^\top \tilde{M}_a \tilde{M}_b \quad p_1^\top \tilde{M}_a \tilde{M}_b \tilde{M}_a \quad \dots$$

How different can they be?

# Belief-continuity is enough

## Definition (Belief-continuity)

A blind MDP is belief-continuous if, for every  $\varepsilon > 0$ , there exists  $\xi > 0$  such that, for all  $\tilde{\Gamma}$  such that  $\text{dist}(\Gamma, \tilde{\Gamma}) \leq \xi$ ,

$$\sup_{\substack{n \geq 1 \\ (a(i))_{i \in [n]}}} \text{dist} \left( M_{a(1)} \cdot \dots \cdot M_{a(n)}, \tilde{M}_{a(1)} \cdot \dots \cdot \tilde{M}_{a(n)} \right) \leq \varepsilon.$$

## Lemma

*If a blind MDP is belief-continuous, then it is XXXX continuous.*

# Belief-continuity

## Theorem

*Every blind MDP is belief continuous.*

Focus on the  $n$ -th step. Define

$$p^\top := p_1^\top M_{a(1)} \cdot \dots \cdot M_{a(n)}$$

$$\tilde{p}^\top := p_1^\top \tilde{M}_{a(1)} \cdot \dots \cdot \tilde{M}_{a(n)}$$

We would like that, for all  $\varepsilon > 0$ , we can choose  $\xi > 0$  so that, for all actions  $a$ ,

$$\text{dist}(p^\top, \tilde{p}^\top) \leq \varepsilon \quad \text{and} \quad \text{dist}(p^\top M_a, \tilde{p}^\top \tilde{M}_a) \leq \varepsilon$$

A stronger notion is the **invariant**

$$\text{dist}(p^\top, \tilde{p}^\top) \leq \varepsilon \quad \Rightarrow \quad \text{dist}(p^\top M_a, \tilde{p}^\top \tilde{M}_a) \leq \varepsilon$$



# Blind MDPs: Belief dynamic

Consider stochastic matrices  $\{M_i : i \in \mathcal{I}\}$   
where the smallest strictly positive transition is uniformly bounded

$$M_i(s, s') > 0 \quad \Rightarrow \quad M_i(s, s') > \delta_{\min} .$$

Consider similar matrices  $\{\tilde{M}_i : i \in \mathcal{I}\}$ .

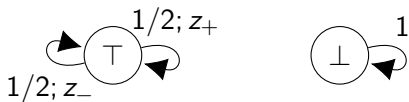
$$p_1^\top \quad p_1^\top M_{i(1)} \quad p_1^\top M_{i(1)} M_{i(2)} \quad p_1^\top M_{i(1)} M_{i(2)} M_{i(3)} \quad \dots$$

$$p_1^\top \quad p_1^\top \tilde{M}_{i(1)} \quad p_1^\top \tilde{M}_{i(1)} \tilde{M}_{i(2)} \quad p_1^\top \tilde{M}_{i(1)} \tilde{M}_{i(2)} \tilde{M}_{i(3)} \quad \dots$$

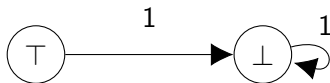
They can not differ by much!

Thank you!

# Motivating example



Action win



Action lose

**Result:** This POMDP is not weakly strategy-continuous.

**Proof:** There is a fragile approximately optimal strategy.

## Proof: Fragile approximately optimal strategy

Consider  $t \geq 1$  large enough and the strategy that plays

$A_1 = A_2 = \dots = A_t = \text{win}$ , and,

if *lose* has been played, then  $A_{m+1} = \text{win}$ ,

if only *win* has been played, for  $m \geq t$ ,

$$A_{m+1} = \text{lose} \quad \Leftrightarrow \quad |\{i \in [2..(m+1)] : Z_i = z_+\}| \geq \left(1 + m^{-1/4}\right) \frac{m}{2}.$$

# Proof: Fragile approximately optimal strategy

## Lemma (Approximate optimality)

Consider  $\Gamma$  the previous POMDP. Then,

$$\mathbb{P}_{p_1}^\sigma[\Gamma](\exists m \geq 1, A_m = \text{lose}) \leq \varepsilon.$$

## Lemma (Fragility)

Consider  $\tilde{\Gamma}$  a small perturbation of  $\Gamma$ . Then,

$$\mathbb{P}_{p_1}^\sigma[\tilde{\Gamma}](\exists m \geq 1, A_m = \text{lose}) = 1.$$

# Extending discontinuity

## Theorem

There exists a POMDP for each of the following combinations.

Example	Continuity		
	Value	Weak strategy	Strong strategy
#1	Yes	Yes	No
#2	No	Yes	No
#3	No	No	No

Remarks:

- All continuities are different
- The exact relationship between the continuity concepts is not fully characterized.

# Characterizing continuity of POMDPs

Theorem (Mathematical characterization, open)

*A POMDP is XXXX continuous if and only if ???*

Theorem (Algorithmic impossibility)

*The problem of deciding whether a given POMDP is XXXX continuous is undecidable.*