

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



K. Chatterjee¹

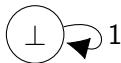


R. Saona¹

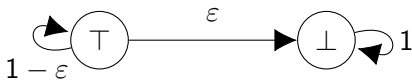
¹Institute of Science and Technology Austria (ISTA)



Continuity in Stochastic dynamics



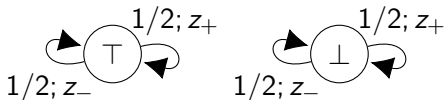
(Deterministic) (dynamic)



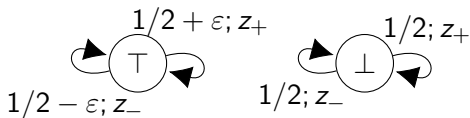
Similar? stochastic dynamic

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	No
POMDPs	No	No	No
Blind MDPs	Yes	Yes	Yes

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 $\Gamma = (\mathcal{S}, \mathcal{A}, \mathcal{Z}, p_1, \delta)$ where

- \mathcal{S} is a finite set of **states**;
- \mathcal{A} is a finite set of **actions**;
- \mathcal{Z} is a finite set of **signals**;
- $p_1 \in \Delta(\mathcal{S})$ is an **initial distribution**;
- $\delta: \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S} \times \mathcal{Z})$ is a probabilistic transition function.

Special cases:

$$|\mathcal{Z}| = 1 \quad \Rightarrow \quad \text{blind MDP}$$

$$\mathcal{Z} = \mathcal{S} \wedge \text{supp}(\delta) \subseteq \{(s, s)\}_{s \in \mathcal{S}} \quad \Rightarrow \quad \text{(fully-observable) MDP}$$

- **strategy** $\sigma: \bigcup_{n \geq 0} (\mathcal{A} \times \mathcal{Z})^n \rightarrow \Delta(\mathcal{A})$
- **play** $\omega = (s_n, a_n, z_{n+1})_{n \geq 1} \subseteq \mathcal{S} \times \mathcal{A} \times \mathcal{Z}$
- **observable history** $h = ((a_i, z_{i+1}))_{i \in [m]} \in (\mathcal{A} \times \mathcal{Z})^m$
- probability $\mathbb{P}_{\rho_1}^\sigma[\Gamma]$ and expectation $\mathbb{E}_{\rho_1}^\sigma[\Gamma]$
- **belief**

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

- **reward** $r: \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$
- **objective** $\gamma(\omega)$ is one of

$$\liminf_{m \rightarrow \infty} \frac{1}{m} \sum_{i=1}^m r(s_i, a_i)$$

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

$$\limsup_{m \rightarrow \infty} \frac{1}{m} \sum_{i=1}^m r(s_i, a_i)$$

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

- set of all strategies \mathcal{X}
- value

$$\text{val}(\Gamma) := \sup_{\sigma \in \mathcal{X}} \mathbb{E}_{p_1}^{\sigma}(\gamma(\omega))$$

- ε -optimal strategy $\mathbb{E}_{p_1}^{\sigma}(\gamma(\omega)) \geq \text{val}(\Gamma) - \varepsilon$ and its set $\mathcal{X}^*(\Gamma, \varepsilon)$
- structural equivalence $\text{supp}(\delta(s, a)) = \text{supp}(\delta'(s, a))$
- ξ -similar POMDPs

$$\sup_{s, a, s', z} |\delta(s, a)(s', z) - \delta'(s, a)(s', z)| \leq \xi$$

Results

Model	Continuity		
	Value	Weak strategy	Strong strategy
Fully-observable MDPs	Yes	Yes	No
POMDPs	No	No	No
Blind MDPs	Yes	Yes	Yes

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

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

Consider an irreducible stochastic matrix Δ .

Computing the stable distribution

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

is a stable operation.

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

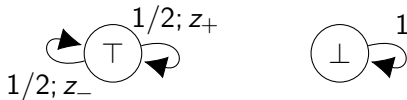
Theorem (Stability of discounted occupation times, Solan 2003)

Consider a Markov Chain with a fixed structure. The λ -discounted occupation time as a function of the transition probabilities is a rational function, i.e., for $\lambda > 0$

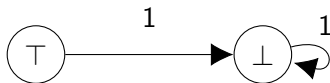
$$\delta \mapsto \text{time}_\lambda(s, \delta) = \frac{\text{poly}(\delta)}{\text{poly}(\delta)}.$$

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

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 Γ the previous POMDP. Then,

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

Lemma (Fragility)

Consider Γ' the previous POMDP. Then,

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

Extending discontinuity

Theorem

There exists a POMDP for each of the following combinations.

<i>Example</i>	<i>Continuity</i>		
	<i>Value</i>	<i>Weak strategy</i>	<i>Strong strategy</i>
<i>#1</i>	<i>Yes</i>	<i>Yes</i>	<i>No</i>
<i>#2</i>	<i>No</i>	<i>Yes</i>	<i>No</i>
<i>#3</i>	<i>No</i>	<i>No</i>	<i>No</i>

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.

Blind MDPs:
no signals guarantee continuity

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 all $\varepsilon > 0$, there exists $\xi > 0$ such that, for all initial belief p_1 , sequence of actions $(a(n))_{n \geq 1}$, and $n \geq 1$

$$\|p_1^\top M_{a(1)} \cdot \dots \cdot M_{a(n)} - p_1^\top \tilde{M}_{a(1)} \cdot \dots \cdot \tilde{M}_{a(n)}\| \leq \varepsilon.$$

Lemma

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

Theorem

Every blind MDP is belief continuous.

Focus on the n -th step. Define

$$\begin{aligned} p^\top &:= p_1^\top M_{a(1)} \cdot \dots \cdot M_{a(n)} \\ q^\top &:= p_1^\top \tilde{M}_{a(1)} \cdot \dots \cdot \tilde{M}_{a(n)} \end{aligned}$$

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

$$\|p^\top - q^\top\| \leq \varepsilon \quad \text{and} \quad \|p^\top M_a - q^\top \tilde{M}_a\| \leq \varepsilon$$

A stronger notion is the **invariant**

$$\|p^\top - q^\top\| \leq \varepsilon \quad \Rightarrow \quad \|p^\top M_a - q^\top \tilde{M}_a\| \leq \varepsilon$$

Lemma

Every blind MDP is belief-continuous as follows.

For every $\varepsilon > 0$, we have that

$$\xi := \varepsilon \frac{\delta_{\min}}{2|\mathcal{S}|}$$

is such that

$$\sup_{\substack{m, h \\ \text{dist}(\Gamma, \Gamma') \leq \xi}} \|P_m[\Gamma](h) - P_m[\Gamma'](h)\|_1 \leq \varepsilon,$$

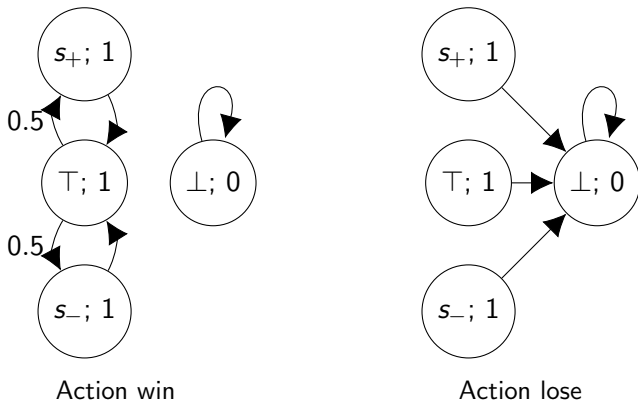
where

$$\delta_{\min} := \min\{\delta(s, a)(s') : a \in \mathcal{A}, s, s' \in \mathcal{S}, \delta(s, a)(s') > 0\},$$

$$\|x\|_1 := \sum_{s \in \mathcal{S}} |x(s)|.$$

Fully-observable MDPs: Fragile ε -optimal strategies

Simulating signals in fully-observable MDPs



There is a fragile approximately-optimal strategy for this MDP.

Thank you!