



Policy Space Identification in Configurable Environments

Alberto Maria Metelli

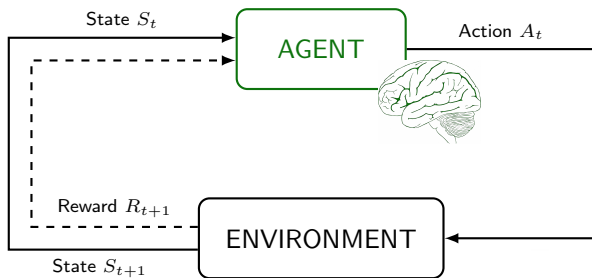
Guglielmo Manneschi

Marcello Restelli

Politecnico di Milano
ECML PKDD 2021 - Journal Track

September 2021

Reinforcement Learning



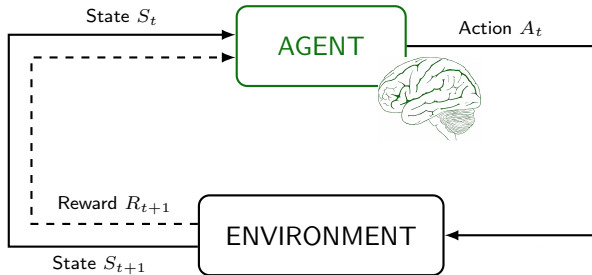
- Markov Decision Process (MDP, Puterman, 2014)
 - 1 Observe the state S_t
 - 2 Perform an action $A_t \sim \pi^{\text{Ag}}(\cdot | S_t)$
 - 3 Transition to the next state $S_{t+1} \sim P(\cdot | S_t, A_t)$
 - 4 Obtain reward $R_{t+1} = r(S_t, A_t, S_{t+1})$

- **Goal:** maximize the expected cumulative discounted reward (Sutton and Barto, 2018):

$$\pi^{\text{Ag}} \in \arg \max_{\pi \in \Pi_{\Theta}} J^{\pi} = \mathbb{E}^{\pi} \left[\sum_{t \in \mathbb{N}} \gamma^t R_{t+1} \right]$$

Π policy space

Reinforcement Learning



- Markov Decision Process (MDP, Puterman, 2014)
 - 1 Observe the state S_t
 - 2 Perform an action $A_t \sim \pi^{\text{Ag}}(\cdot | S_t)$
 - 3 Transition to the next state $S_{t+1} \sim P(\cdot | S_t, A_t)$
 - 4 Obtain reward $R_{t+1} = r(S_t, A_t, S_{t+1})$

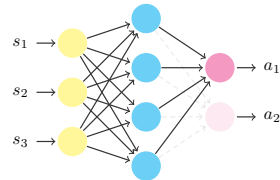
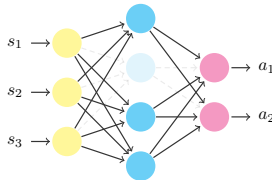
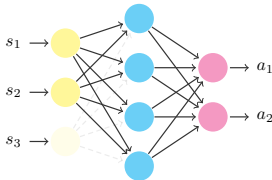
- **Goal:** maximize the expected cumulative discounted reward (Sutton and Barto, 2018):

$$\pi^{\text{Ag}} \in \arg \max_{\pi \in \Pi_{\Theta}} J^{\pi} = \mathbb{E}^{\pi} \left[\sum_{t \in \mathbb{N}} \gamma^t R_{t+1} \right]$$

Π policy space

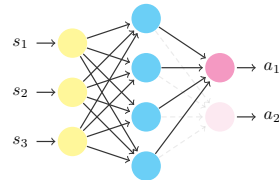
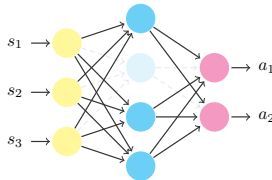
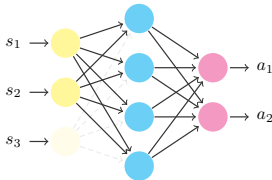
Motivations and Problem

- The **policy space** Π defines the **perception**, **actuation**, and **mapping** capabilities of an agent
- **Research Question:** How to identify the **policy space** of an agent by observing its behavior $\pi^{\text{Ag?}}$? → **Policy Space Identification (PSI)**
- Applications
 - Configurable MDPs (Metelli et al., 2018)
 - Imitation Learning (Osa et al., 2018)



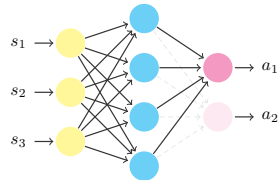
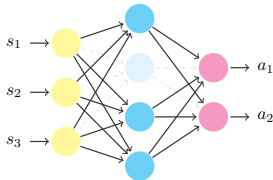
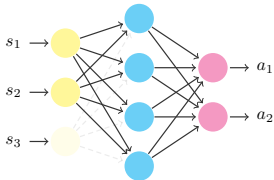
Motivations and Problem

- The **policy space** Π defines the **perception**, **actuation**, and **mapping** capabilities of an agent
- **Research Question:** How to identify the **policy space** of an agent by observing its behavior $\pi^{\text{Ag?}}$ → **Policy Space Identification (PSI)**
- Applications
 - Configurable MDPs (Metelli et al., 2018)
 - Imitation Learning (Osa et al., 2018)



Motivations and Problem

- The **policy space** Π defines the **perception**, **actuation**, and **mapping** capabilities of an agent
- **Research Question:** How to identify the **policy space** of an agent by observing its behavior $\pi^{\text{Ag?}}$ → **Policy Space Identification (PSI)**
- Applications
 - Configurable MDPs (Metelli et al., 2018)
 - Imitation Learning (Osa et al., 2018)



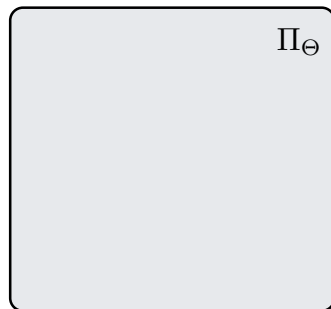
Policy Spaces and Correctness

- **Parametric** policy space Π_{Θ} , with $\Theta \subset \mathbb{R}^d$
- Agent's policy $\pi^{\text{Ag}} \in \Pi_{\Theta}$
- The agent **controls** (i.e., can change) $d^{\text{Ag}} < d$ parameters
- $I \subseteq \{1, \dots, d\}$ subset of indexes

$$\Theta_I = \{\theta \in \Theta : \theta_i = 0, \forall i \in \{1, \dots, d\} \setminus I\}$$

- I^{Ag} is **correct** for the agent's policy π^{Ag} iff

$$\underbrace{\pi^{\text{Ag}} \in \Pi_{\Theta_{I^{\text{Ag}}}}}_{\text{sufficient}} \wedge \underbrace{\forall i \in I^{\text{Ag}} : \pi^{\text{Ag}} \notin \Pi_{\Theta_{I^{\text{Ag}} \setminus \{i\}}}}_{\text{necessary}}$$



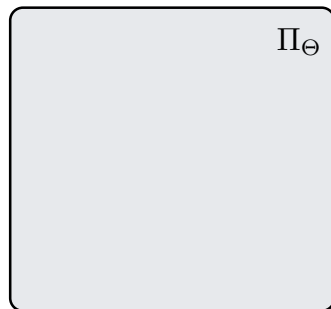
Policy Spaces and Correctness

- **Parametric** policy space Π_{Θ} , with $\Theta \subset \mathbb{R}^d$
- Agent's policy $\pi^{\text{Ag}} \in \Pi_{\Theta}$
- The agent **controls** (i.e., can change) $d^{\text{Ag}} < d$ parameters
- $I \subseteq \{1, \dots, d\}$ subset of indexes

$$\Theta_I = \{\theta \in \Theta : \theta_i = 0, \forall i \in \{1, \dots, d\} \setminus I\}$$

- I^{Ag} is **correct** for the agent's policy π^{Ag} iff

$$\underbrace{\pi^{\text{Ag}} \in \Pi_{\Theta_{I^{\text{Ag}}}}}_{\text{sufficient}} \wedge \underbrace{\forall i \in I^{\text{Ag}} : \pi^{\text{Ag}} \notin \Pi_{\Theta_{I^{\text{Ag}} \setminus \{i\}}}}_{\text{necessary}}$$



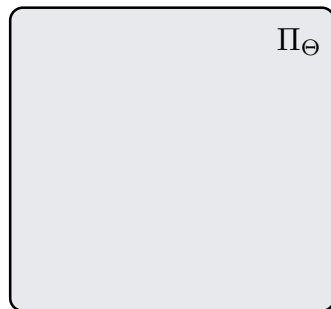
Policy Spaces and Correctness

- **Parametric** policy space Π_{Θ} , with $\Theta \subset \mathbb{R}^d$
- Agent's policy $\pi^{\text{Ag}} \in \Pi_{\Theta}$
- The agent **controls** (i.e., can change) $d^{\text{Ag}} < d$ parameters
- $I \subseteq \{1, \dots, d\}$ subset of indexes

$$\Theta_I = \{\theta \in \Theta : \theta_i = 0, \forall i \in \{1, \dots, d\} \setminus I\}$$

- I^{Ag} is **correct** for the agent's policy π^{Ag} iff

$$\underbrace{\pi^{\text{Ag}} \in \Pi_{\Theta_{I^{\text{Ag}}}}}_{\text{sufficient}} \wedge \underbrace{\forall i \in I^{\text{Ag}} : \pi^{\text{Ag}} \notin \Pi_{\Theta_{I^{\text{Ag}} \setminus \{i\}}}}_{\text{necessary}}$$



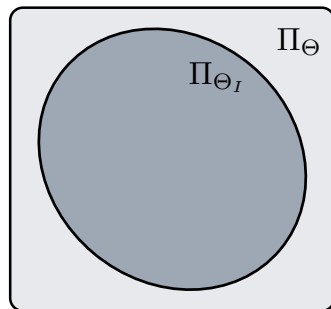
Policy Spaces and Correctness

- **Parametric** policy space Π_{Θ} , with $\Theta \subset \mathbb{R}^d$
- Agent's policy $\pi^{\text{Ag}} \in \Pi_{\Theta}$
- The agent **controls** (i.e., can change) $d^{\text{Ag}} < d$ parameters
- $I \subseteq \{1, \dots, d\}$ subset of indexes

$$\Theta_I = \{\theta \in \Theta : \theta_i = 0, \forall i \in \{1, \dots, d\} \setminus I\}$$

- I^{Ag} is **correct** for the agent's policy π^{Ag} iff

$$\underbrace{\pi^{\text{Ag}} \in \Pi_{\Theta_{I^{\text{Ag}}}}}_{\text{sufficient}} \wedge \underbrace{\forall i \in I^{\text{Ag}} : \pi^{\text{Ag}} \notin \Pi_{\Theta_{I^{\text{Ag}} \setminus \{i\}}}}_{\text{necessary}}$$



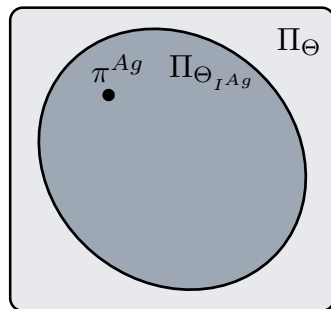
Policy Spaces and Correctness

- **Parametric** policy space Π_{Θ} , with $\Theta \subset \mathbb{R}^d$
- Agent's policy $\pi^{\text{Ag}} \in \Pi_{\Theta}$
- The agent **controls** (i.e., can change) $d^{\text{Ag}} < d$ parameters
- $I \subseteq \{1, \dots, d\}$ subset of indexes

$$\Theta_I = \{\theta \in \Theta : \theta_i = 0, \forall i \in \{1, \dots, d\} \setminus I\}$$

- I^{Ag} is **correct** for the agent's policy π^{Ag} iff

$$\underbrace{\pi^{\text{Ag}} \in \Pi_{\Theta_{I^{\text{Ag}}}}}_{\text{sufficient}} \wedge \underbrace{\forall i \in I^{\text{Ag}} : \pi^{\text{Ag}} \notin \Pi_{\Theta_{I^{\text{Ag}} \setminus \{i\}}}}_{\text{necessary}}$$



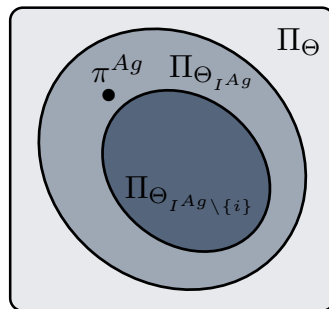
Policy Spaces and Correctness

- **Parametric** policy space Π_{Θ} , with $\Theta \subset \mathbb{R}^d$
- Agent's policy $\pi^{\text{Ag}} \in \Pi_{\Theta}$
- The agent **controls** (i.e., can change) $d^{\text{Ag}} < d$ parameters
- $I \subseteq \{1, \dots, d\}$ subset of indexes

$$\Theta_I = \{\theta \in \Theta : \theta_i = 0, \forall i \in \{1, \dots, d\} \setminus I\}$$

- I^{Ag} is **correct** for the agent's policy π^{Ag} iff

$$\underbrace{\pi^{\text{Ag}} \in \Pi_{\Theta_{I^{\text{Ag}}}}}_{\text{sufficient}} \quad \wedge \quad \underbrace{\forall i \in I^{\text{Ag}} : \pi^{\text{Ag}} \notin \Pi_{\Theta_{I^{\text{Ag}} \setminus \{i\}}}}_{\text{necessary}}$$



Hypothesis Tests

- **Idea:** perform **hypothesis test** for $I \subseteq \{1, \dots, d\}$

$$\mathcal{H}_{0,I} : \pi^{\text{Ag}} \in \Pi_{\Theta_I} \quad \text{vs} \quad \mathcal{H}_{1,I} : \pi^{\text{Ag}} \in \Pi_{\Theta \setminus \Theta_I}$$

- Dataset of samples $\{(S_i, A_i)\}_{i=1}^n$ collected with the agent's policy π^{Ag}
- Likelihood of a parameter $\theta \in \Theta$

$$\hat{\mathcal{L}}(\theta) = \prod_{i=1}^n \pi_{\theta}(A_i | S_i) \qquad \mathcal{L}(\theta) = \mathbb{E}[\hat{\mathcal{L}}(\theta)]$$

- Generalized **likelihood ratio** statistic (Casella and Berger, 2002)

$$\Lambda_I = \frac{\sup_{\theta \in \Theta_I} \hat{\mathcal{L}}(\theta)}{\sup_{\theta \in \Theta} \hat{\mathcal{L}}(\theta)}$$

$$\Lambda_I \simeq 0 \rightarrow \text{reject } \mathcal{H}_{0,I}$$

$$\Lambda_I \simeq 1 \rightarrow \text{do not reject } \mathcal{H}_{0,I}$$

Hypothesis Tests

- **Idea:** perform **hypothesis test** for $I \subseteq \{1, \dots, d\}$

$$\mathcal{H}_{0,I} : \pi^{\text{Ag}} \in \Pi_{\Theta_I} \quad \text{vs} \quad \mathcal{H}_{1,I} : \pi^{\text{Ag}} \in \Pi_{\Theta \setminus \Theta_I}$$

- Dataset of samples $\{(S_i, A_i)\}_{i=1}^n$ collected with the agent's policy π^{Ag}
- Likelihood of a parameter $\theta \in \Theta$

$$\hat{\mathcal{L}}(\theta) = \prod_{i=1}^n \pi_{\theta}(A_i | S_i) \qquad \mathcal{L}(\theta) = \mathbb{E}[\hat{\mathcal{L}}(\theta)]$$

- Generalized **likelihood ratio** statistic (Casella and Berger, 2002)

$$\Lambda_I = \frac{\sup_{\theta \in \Theta_I} \hat{\mathcal{L}}(\theta)}{\sup_{\theta \in \Theta} \hat{\mathcal{L}}(\theta)}$$

$$\Lambda_I \simeq 0 \rightarrow \text{reject } \mathcal{H}_{0,I}$$

$$\Lambda_I \simeq 1 \rightarrow \text{do not reject } \mathcal{H}_{0,I}$$

Hypothesis Tests

- **Idea:** perform **hypothesis test** for $I \subseteq \{1, \dots, d\}$

$$\mathcal{H}_{0,I} : \pi^{\text{Ag}} \in \Pi_{\Theta_I} \quad \text{vs} \quad \mathcal{H}_{1,I} : \pi^{\text{Ag}} \in \Pi_{\Theta \setminus \Theta_I}$$

- Dataset of samples $\{(S_i, A_i)\}_{i=1}^n$ collected with the agent's policy π^{Ag}
- Likelihood of a parameter $\theta \in \Theta$

$$\hat{\mathcal{L}}(\theta) = \prod_{i=1}^n \pi_{\theta}(A_i | S_i) \qquad \mathcal{L}(\theta) = \mathbb{E}[\hat{\mathcal{L}}(\theta)]$$

- Generalized **likelihood ratio** statistic (Casella and Berger, 2002)

$$\Lambda_I = \frac{\sup_{\theta \in \Theta_I} \hat{\mathcal{L}}(\theta)}{\sup_{\theta \in \Theta} \hat{\mathcal{L}}(\theta)}$$

$$\Lambda_I \simeq 0 \rightarrow \text{reject } \mathcal{H}_{0,I}$$

$$\Lambda_I \simeq 1 \rightarrow \text{do not reject } \mathcal{H}_{0,I}$$

Hypothesis Tests

- **Idea:** perform **hypothesis test** for $I \subseteq \{1, \dots, d\}$

$$\mathcal{H}_{0,I} : \pi^{\text{Ag}} \in \Pi_{\Theta_I} \quad \text{vs} \quad \mathcal{H}_{1,I} : \pi^{\text{Ag}} \in \Pi_{\Theta \setminus \Theta_I}$$

- Dataset of samples $\{(S_i, A_i)\}_{i=1}^n$ collected with the agent's policy π^{Ag}
- Likelihood of a parameter $\theta \in \Theta$

$$\hat{\mathcal{L}}(\theta) = \prod_{i=1}^n \pi_{\theta}(A_i | S_i) \qquad \mathcal{L}(\theta) = \mathbb{E}[\hat{\mathcal{L}}(\theta)]$$

- Generalized **likelihood ratio** statistic (Casella and Berger, 2002)

$$\Lambda_I = \frac{\sup_{\theta \in \Theta_I} \hat{\mathcal{L}}(\theta)}{\sup_{\theta \in \Theta} \hat{\mathcal{L}}(\theta)}$$

$$\Lambda_I \simeq 0 \rightarrow \text{reject } \mathcal{H}_{0,I}$$

$$\Lambda_I \simeq 1 \rightarrow \text{do not reject } \mathcal{H}_{0,I}$$

Hypothesis Tests

- **Idea:** perform **hypothesis test** for $I \subseteq \{1, \dots, d\}$

$$\mathcal{H}_{0,I} : \pi^{\text{Ag}} \in \Pi_{\Theta_I} \quad \text{vs} \quad \mathcal{H}_{1,I} : \pi^{\text{Ag}} \in \Pi_{\Theta \setminus \Theta_I}$$

- Dataset of samples $\{(S_i, A_i)\}_{i=1}^n$ collected with the agent's policy π^{Ag}
- Likelihood of a parameter $\theta \in \Theta$

$$\hat{\mathcal{L}}(\theta) = \prod_{i=1}^n \pi_{\theta}(A_i | S_i) \qquad \mathcal{L}(\theta) = \mathbb{E}[\hat{\mathcal{L}}(\theta)]$$

- Generalized **likelihood ratio** statistic (Casella and Berger, 2002)

$$\Lambda_I = \frac{\sup_{\theta \in \Theta_I} \hat{\mathcal{L}}(\theta)}{\sup_{\theta \in \Theta} \hat{\mathcal{L}}(\theta)}$$

$$\Lambda_I \simeq 0 \rightarrow \text{reject } \mathcal{H}_{0,I}$$

$$\Lambda_I \simeq 1 \rightarrow \text{do not reject } \mathcal{H}_{0,I}$$

Combinatorial Identification Rule

- **Combinatorial Identification Rule:** retain all the **approximately correct** $\hat{I} \subseteq \{1, \dots, d\}$:

$$\underbrace{\text{do not reject } \mathcal{H}_{0,\hat{I}}}_{\text{sufficient}} \quad \wedge \quad \underbrace{\forall i \in \hat{I} : \text{reject } \mathcal{H}_{0,\hat{I} \setminus \{i\}}}_{\text{necessary}}$$

- Requires $O(2^d)$ tests! \rightarrow combinatorial
- Works under multiple representations of π^{Ag}

What if there exists a unique representation of π^{Ag} ?

Combinatorial Identification Rule

- **Combinatorial Identification Rule:** retain all the **approximately correct** $\hat{I} \subseteq \{1, \dots, d\}$:

$$\underbrace{\text{do not reject } \mathcal{H}_{0,\hat{I}}}_{\text{sufficient}} \quad \wedge \quad \underbrace{\forall i \in \hat{I} : \text{reject } \mathcal{H}_{0,\hat{I} \setminus \{i\}}}_{\text{necessary}}$$

- Requires $O(2^d)$ tests! \rightarrow combinatorial
- Works under multiple representations of π^{Ag}

What if there exists a unique representation of π^{Ag} ?

Combinatorial Identification Rule

- **Combinatorial Identification Rule:** retain all the **approximately correct** $\hat{I} \subseteq \{1, \dots, d\}$:

$$\underbrace{\text{do not reject } \mathcal{H}_{0, \hat{I}}}_{\text{sufficient}} \quad \wedge \quad \underbrace{\forall i \in \hat{I} : \text{reject } \mathcal{H}_{0, \hat{I} \setminus \{i\}}}_{\text{necessary}}$$

- Requires $O(2^d)$ tests! \rightarrow combinatorial
- Works under multiple representations of π^{Ag}

What if there exists a unique representation of π^{Ag} ?

Combinatorial Identification Rule

- **Combinatorial Identification Rule:** retain all the **approximately correct** $\hat{I} \subseteq \{1, \dots, d\}$:

$$\underbrace{\text{do not reject } \mathcal{H}_{0,\hat{I}}}_{\text{sufficient}} \quad \wedge \quad \underbrace{\forall i \in \hat{I} : \text{reject } \mathcal{H}_{0,\hat{I} \setminus \{i\}}}_{\text{necessary}}$$

- Requires $O(2^d)$ tests! \rightarrow combinatorial
- Works under multiple representations of π^{Ag}

What if there exists a unique representation of π^{Ag} ?

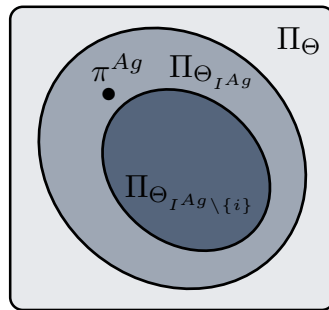
Identifiability and Correctness

- Π_{Θ} is **identifiable** if

$$\forall \theta', \theta \in \Theta \quad \pi_{\theta'}(\cdot|s) = \pi_{\theta}(\cdot|s) \text{ a.s.} \implies \theta' = \theta$$

- We can reason on the **parameters** only!
- The **only** correct I^{Ag} for the agent's policy $\pi_{\theta^{\text{Ag}}}$ is

$$I^{\text{Ag}} = \left\{ i \in \{1, \dots, d\} : \theta_i^{\text{Ag}} \neq 0 \right\}$$



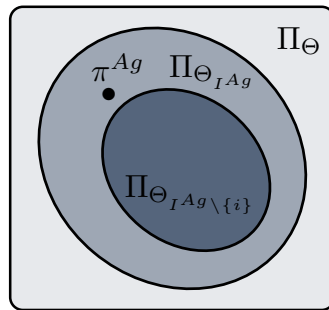
Identifiability and Correctness

- Π_{Θ} is **identifiable** if

$$\forall \theta', \theta \in \Theta \quad \pi_{\theta'}(\cdot|s) = \pi_{\theta}(\cdot|s) \text{ a.s.} \implies \theta' = \theta$$

- We can reason on the **parameters** only!
- The **only** correct I^{Ag} for the agent's policy $\pi_{\theta^{\text{Ag}}}$ is

$$I^{\text{Ag}} = \left\{ i \in \{1, \dots, d\} : \theta_i^{\text{Ag}} \neq 0 \right\}$$



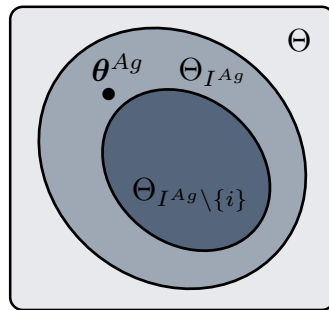
Identifiability and Correctness

- Π_{Θ} is **identifiable** if

$$\forall \theta', \theta \in \Theta \quad \pi_{\theta'}(\cdot|s) = \pi_{\theta}(\cdot|s) \text{ a.s.} \implies \theta' = \theta$$

- We can reason on the **parameters** only!
- The **only** correct I^{Ag} for the agent's policy $\pi_{\theta^{\text{Ag}}}$ is

$$I^{\text{Ag}} = \left\{ i \in \{1, \dots, d\} : \theta_i^{\text{Ag}} \neq 0 \right\}$$



Simplified Identification Rule

- **Idea:** perform **hypothesis test** for $i \in \{1, \dots, d\}$

$$\mathcal{H}_{0,i} : \theta_i^{\text{Ag}} = 0 \quad \text{vs} \quad \mathcal{H}_{1,i} : \theta_i^{\text{Ag}} \neq 0$$

- **Simplified Identification Rule:** \hat{I} is the union of all $i \in \{1, \dots, d\}$ such that:

reject $\mathcal{H}_{0,i}$

- Requires $O(d)$ tests! \rightarrow linear
- Works under **identifiability** only!

Simplified Identification Rule

- **Idea:** perform **hypothesis test** for $i \in \{1, \dots, d\}$

$$\mathcal{H}_{0,i} : \theta_i^{\text{Ag}} = 0 \quad \text{vs} \quad \mathcal{H}_{1,i} : \theta_i^{\text{Ag}} \neq 0$$

- **Simplified Identification Rule:** \hat{I} is the union of all $i \in \{1, \dots, d\}$ such that:

reject $\mathcal{H}_{0,i}$

- Requires $O(d)$ tests! \rightarrow linear
- Works under **identifiability** only!

Simplified Identification Rule

- **Idea:** perform **hypothesis test** for $i \in \{1, \dots, d\}$

$$\mathcal{H}_{0,i} : \theta_i^{\text{Ag}} = 0 \quad \text{vs} \quad \mathcal{H}_{1,i} : \theta_i^{\text{Ag}} \neq 0$$

- **Simplified Identification Rule:** \hat{I} is the union of all $i \in \{1, \dots, d\}$ such that:

reject $\mathcal{H}_{0,i}$

- Requires $O(d)$ tests! \rightarrow linear
- Works under **identifiability** only!

Simplified Identification Rule

- **Idea:** perform **hypothesis test** for $i \in \{1, \dots, d\}$

$$\mathcal{H}_{0,i} : \theta_i^{\text{Ag}} = 0 \quad \text{vs} \quad \mathcal{H}_{1,i} : \theta_i^{\text{Ag}} \neq 0$$

- **Simplified Identification Rule:** \hat{I} is the union of all $i \in \{1, \dots, d\}$ such that:

reject $\mathcal{H}_{0,i}$

- Requires $O(d)$ tests! \rightarrow linear
- Works under **identifiability** only!

Ambiguous Identification

What can we conclude when $\theta_i = 0$?

- ① The agent **does not control** θ_i or...
 - ② ...the agent has **consciously** chosen to set $\theta_i = 0$
- **Problem:** How to distinguish between these two scenarios?
 - **Idea:** change the environment to make the parameter “maximally important” → **Configurable Environment** (Metelli et al., 2018)

Ambiguous Identification

What can we conclude when $\theta_i = 0$?

- ① The agent **does not control** θ_i or...
 - ② ...the agent has **consciously** chosen to set $\theta_i = 0$
- **Problem:** How to distinguish between these two scenarios?
 - **Idea:** change the environment to make the parameter “maximally important” → **Configurable Environment** (Metelli et al., 2018)

Ambiguous Identification

What can we conclude when $\theta_i = 0$?

- ① The agent **does not control** θ_i or...
 - ② ...the agent has **consciously** chosen to set $\theta_i = 0$
- **Problem:** How to distinguish between these two scenarios?
 - **Idea:** change the environment to make the parameter “maximally important” → **Configurable Environment** (Metelli et al., 2018)

Ambiguous Identification

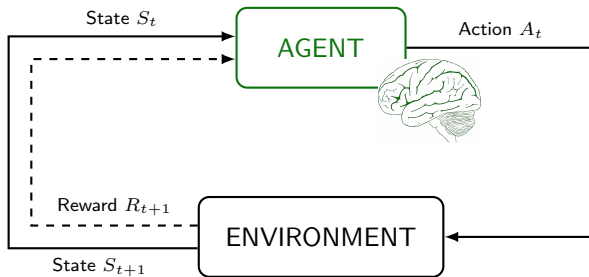
What can we conclude when $\theta_i = 0$?

- ① The agent **does not control** θ_i or...
 - ② ...the agent has **consciously** chosen to set $\theta_i = 0$
-
- **Problem:** How to distinguish between these two scenarios?
 - **Idea:** change the environment to make the parameter “maximally important” → **Configurable Environment** (Metelli et al., 2018)

Non-Configurable Environments

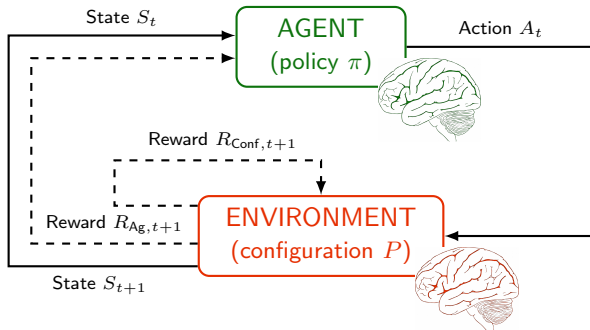
MDP (Puterman, 2014)

the environment is fixed and out of control



Configurable Environments

Conf-MDP (Metelli et al., 2018)
the environment can be configured



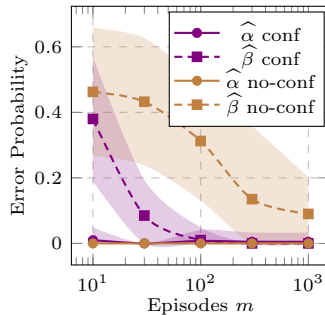
Performance of the Identification Rules

Discrete Grid world

- π_θ linear in RBF features ϕ
- Agent observes a limited subset of ϕ
- Configure the initial state distribution

$$\alpha = \Pr\left(\exists i \notin I^{\text{Ag}} : i \in \hat{I}\right)$$

$$\beta = \Pr\left(\exists i \in I^{\text{Ag}} : i \notin \hat{I}\right)$$

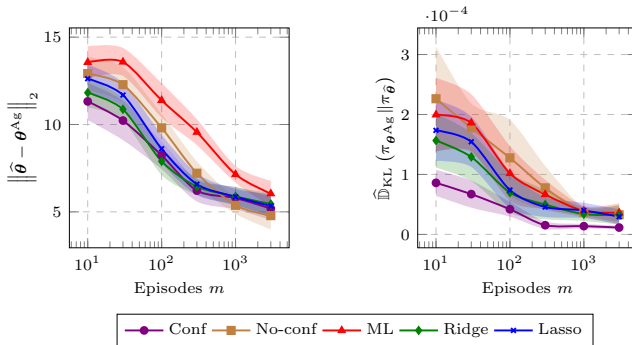


Application to Imitation Learning

Discrete Grid world

- Maximum likelihood policy estimation

$$-\sum_{i=1}^n \log \pi_{\theta}(a_i | s_i) + \text{Reg}(\theta)$$

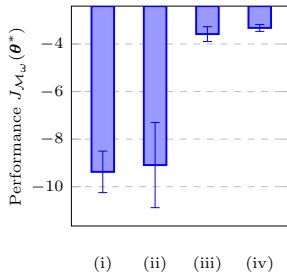
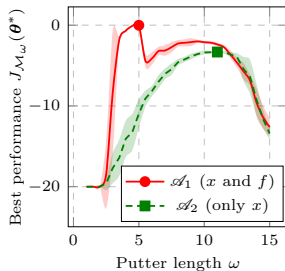


Application to Conf-MDP

Minigolf

- \mathcal{A}_1 observes *position* x and *friction* f
- \mathcal{A}_2 observes *position* x only

- (i) random choice
- (ii) optimal for \mathcal{A}_1
- (iii) using identification rule
- (iv) optimal for \mathcal{A}_2



Discussion and Conclusions

Contributions

- Identification rules
- Environment configurability to improve identification
- Applications of policy space identification

Future Works

- Bayesian statistical tests
- Multi-agent systems

Discussion and Conclusions

Contributions

- Identification rules
- Environment configurability to improve identification
- Applications of policy space identification

Future Works

- Bayesian statistical tests
- Multi-agent systems

Thank You for Your Attention!

Code: `github.com/albertometelli/policy-space-identification`

Contact: `albertomaria.metelli@polimi.it`

References I

George Casella and Roger L Berger. *Statistical inference*, volume 2. Duxbury Pacific Grove, CA, 2002.

Alberto Maria Metelli, Mirco Mutti, and Marcello Restelli. Configurable markov decision processes. In Jennifer G. Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pages 3488–3497. PMLR, 2018.

Takayuki Osa, Joni Pajarinen, Gerhard Neumann, J. Andrew Bagnell, Pieter Abbeel, and Jan Peters. An algorithmic perspective on imitation learning. *Foundations and Trends in Robotics*, 7(1-2):1–179, 2018. doi: 10.1561/23000000053.

Martin L Puterman. *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons, 2014.

Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.