



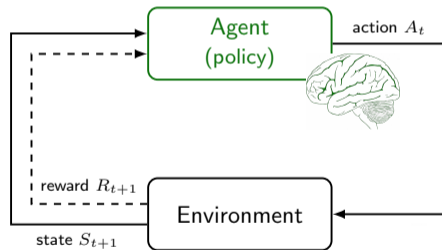
POLITECNICO
MILANO 1863

Reinforcement Learning in Configurable Continuous Environments

Alberto Maria Metelli, Emanuele Ghelfi and Marcello Restelli

36th International Conference on Machine Learning
13th June 2019

Markov Decision Process
(MDP, Puterman, 2014)



$$\mathcal{M} = (\mathcal{S}, \mathcal{A}, r, \gamma, \mu, p)$$

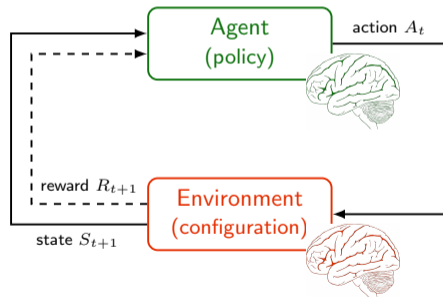
$$S_0 \sim \mu, A_t \sim \pi_{\theta}(\cdot | S_t), S_{t+1} \sim p(\cdot | S_t, A_t)$$

- Learn the **policy parameters** θ under the fixed environment p

$$\theta^* = \arg \max_{\theta \in \Theta} J(\theta) = \mathbb{E} \left[\sum_{t=0}^{+\infty} \gamma^t R_{t+1} \right]$$



Configurable Markov Decision Process (Conf-MDP, Metelli et al., 2018)



$$\mathcal{CM} = (\mathcal{S}, \mathcal{A}, r, \gamma, \mu, \mathcal{P}, \Pi)$$

$$S_0 \sim \mu, A_t \sim \pi_{\theta}(\cdot | S_t), S_{t+1} \sim p_{\omega}(\cdot | S_t, A_t)$$

- Learn the **policy parameters** θ together with the **environment configuration** ω

$$\theta^*, \omega^* = \arg \max_{\theta \in \Theta, \omega \in \Omega} J(\theta, \omega) = \mathbb{E} \left[\sum_{t=0}^{+\infty} \gamma^t R_{t+1} \right]$$

- **Safe Policy Model Iteration** (SPMI, Metelli et al., 2018)
 - Optimize a lower bound of the performance improvement
- Limitations
 - **Finite** state-actions spaces
 - **Full knowledge** of the environment dynamics
- Similar approaches Keren et al. (2017) and Silva et al. (2018)

Optimization

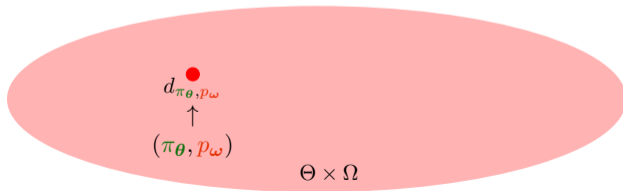
Find a new *stationary distribution* d' in a *trust region* centered in $d_{\pi_{\theta}, p_{\omega}}$

$$\begin{aligned} \max_{d'} J_{d'} &= \mathbb{E}_{S, A, S' \sim d'} [r(S, A, S')] \\ \text{s.t. } D_{\text{KL}}(d' \| d_{\pi_{\theta}, p_{\omega}}) &\leq \kappa, \end{aligned}$$

Projection

Find a policy $\pi_{\theta'}$ and configuration $p_{\omega'}$ inducing a stationary distribution close to d'

$$\min_{\theta' \in \Theta, \omega' \in \Omega} D_{\text{KL}}(d' \| d_{\pi_{\theta'}, p_{\omega'}})$$



Optimization

Find a new *stationary distribution* d' in a *trust region* centered in $d_{\pi_{\theta}, p_{\omega}}$

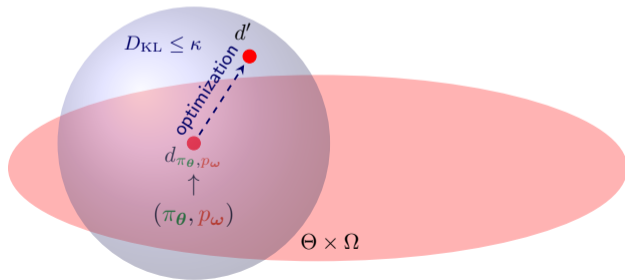
$$\max_{d'} J_{d'} = \mathbb{E}_{S, A, S' \sim d'} [r(S, A, S')]$$

$$\text{s.t. } D_{\text{KL}}(d' \| d_{\pi_{\theta}, p_{\omega}}) \leq \kappa,$$

Projection

Find a policy $\pi_{\theta'}$ and configuration $p_{\omega'}$ inducing a stationary distribution close to d'

$$\min_{\theta' \in \Theta, \omega' \in \Omega} D_{\text{KL}}(d' \| d_{\pi_{\theta'}, p_{\omega'}})$$



Optimization

Find a new *stationary distribution* d' in a *trust region* centered in $d_{\pi_{\theta}, p_{\omega}}$

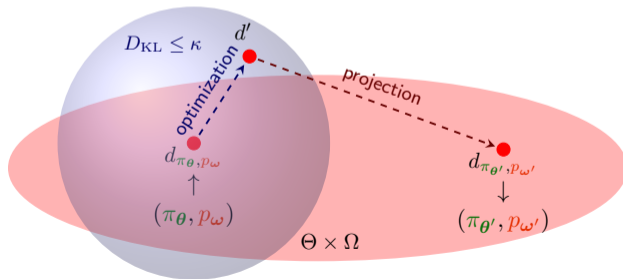
$$\max_{d'} J_{d'} = \mathbb{E}_{S, A, S' \sim d'} [r(S, A, S')]$$

$$\text{s.t. } D_{\text{KL}}(d' \| d_{\pi_{\theta}, p_{\omega}}) \leq \kappa,$$

Projection

Find a policy $\pi_{\theta'}$ and configuration $p_{\omega'}$ inducing a stationary distribution close to d'

$$\min_{\theta' \in \Theta, \omega' \in \Omega} D_{\text{KL}}(d' \| d_{\pi_{\theta'}, p_{\omega'}})$$



Optimization

Find a new *stationary distribution* d' in a *trust region* centered in $d_{\pi_{\theta}, p_{\omega}}$

$$\max_{d'} J_{d'} = \mathbb{E}_{S, A, S' \sim d'} [r(S, A, S')]$$

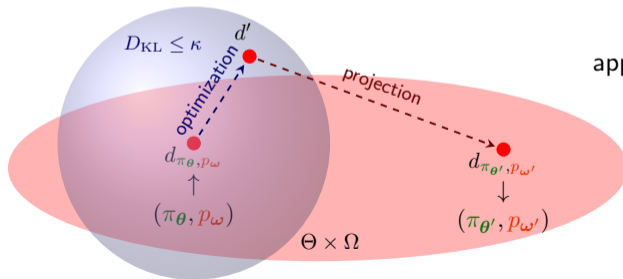
$$\text{s.t. } D_{\text{KL}}(d' \| d_{\pi_{\theta}, p_{\omega}}) \leq \kappa,$$

Projection

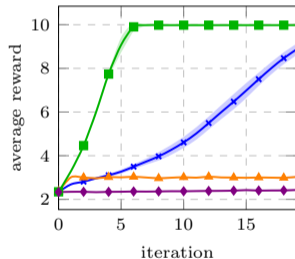
Find a policy $\pi_{\theta'}$ and configuration $p_{\omega'}$ inducing a stationary distribution close to d'

$$\min_{\theta' \in \Theta, \omega' \in \Omega} D_{\text{KL}}(d' \| d_{\pi_{\theta'}, p_{\omega'}})$$

Can also be an approximated model \hat{p}



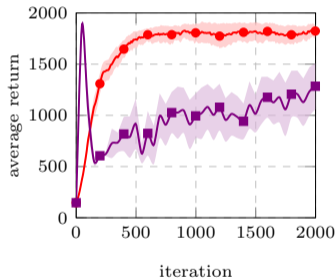
Chain Domain



—●— REMPS (0.01) —■— REMPS (0.1)
—▲— REMPS (10) —◆— G(PO)MDP

Cartpole

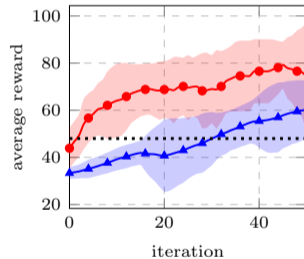
■ Configure the *cart force*



—●— REMPS —■— G(PO)MDP

TORCS

■ Configure the *front-rear wing orientation* and *brake repartition*



—●— REMPS —▲— REPS
 Bot

Thank You for Your Attention!

- Poster **Pacific Ballroom #37**
- Code: `github.com/albertometelli/remps`
- Web page: `albertometelli.github.io/ICML2019-REMPS`
- Contact: `albertomaria.metelli@polimi.it`



- Keren, S., Pineda, L., Gal, A., Karpas, E., and Zilberstein, S. (2017). Equi-reward utility maximizing design in stochastic environments. In *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI-17*, pages 4353–4360.
- Metelli, A. M., Mutti, M., and Restelli, M. (2018). Configurable markov decision processes. In Dy, J. G. and Krause, A., 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.
- Puterman, M. L. (2014). *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons.
- Silva, R., Melo, F. S., and Veloso, M. (2018). What if the world were different? gradient-based exploration for new optimal policies. *EPiC Series in Computing*, 55:229–242.