



POLITECNICO
MILANO 1863

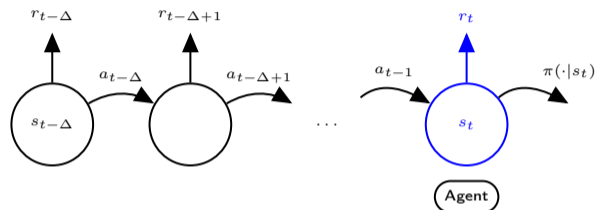
Delayed Reinforcement Learning by Imitation

Pierre Liotet Davide Maran Lorenzo Bisi Marcello Restelli

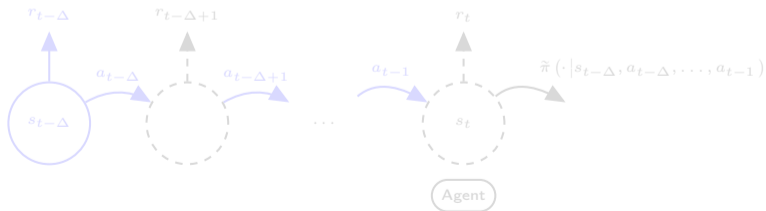
July 2022

Thirty-ninth International Conference on Machine Learning (ICML-22)

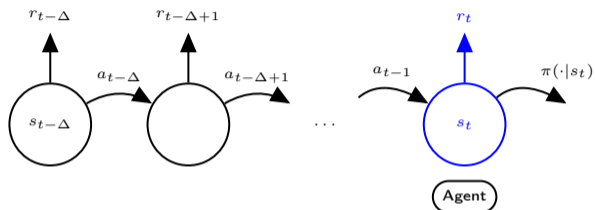
■ Reinforcement learning:



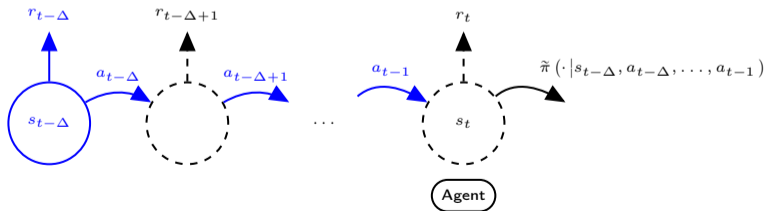
■ Delay:



■ Reinforcement learning:



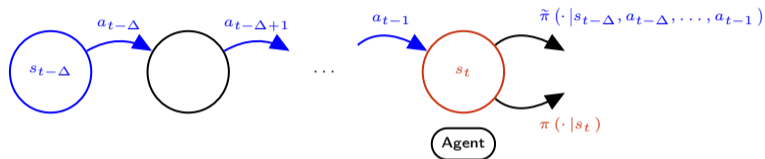
■ Delay:



■ Complications:

- Non-Markovianity
- Augmented state space
- Performance loss
- Non integer delays

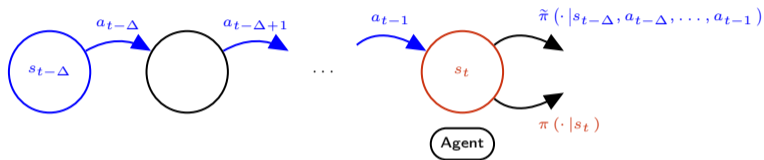
- Duality of trajectories



- Imitation learning with DAgger (Ross et al., 2011):

$$\tilde{\pi}(a | s_{t-\Delta}, a_{t-\Delta}, \dots, a_{t-1}) = \int_{\mathcal{S}} b(s | s_{t-\Delta}, a_{t-\Delta}, \dots, a_{t-1}) \pi(a | s) ds.$$

- Duality of trajectories



- Imitation learning with DAgger (Ross et al., 2011):

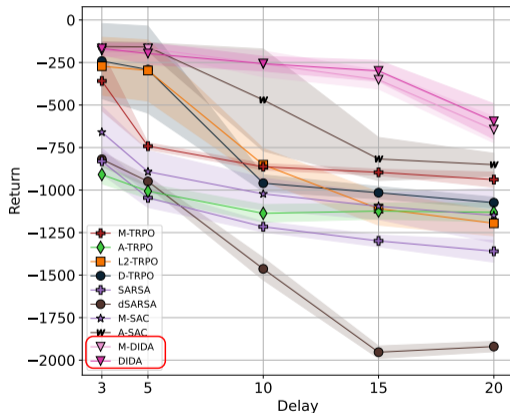
$$\tilde{\pi}(a | s_{t-\Delta}, a_{t-\Delta}, \dots, a_{t-1}) = \int_{\mathcal{S}} b(s | s_{t-\Delta}, a_{t-\Delta}, \dots, a_{t-1}) \pi(a | s) ds.$$

- Bounding performance loss:

$$\underbrace{\mathbb{E}_{s \sim b(\cdot|x)} [V^\pi(s)]}_{\text{undelayed expert performance}} - \underbrace{V^{\tilde{\pi}}(x)}_{\text{delayed performance}} \leq \frac{C}{1-\gamma} \underbrace{\mathbb{E}_{x' \sim d_x^{\tilde{\pi}}(\cdot)} \left[\sqrt{\text{Var}_{s \sim b(\cdot|x')} (s|x')} \right]}_{\text{measure of uncertainty on the current unobserved state}}.$$

- Deterministic process \implies bound = 0

■ Mean return as a function delay



- Robotic locomotion - mujoco (Todorov et al., 2012)
- Trading

- S. Ross, G. Gordon, and D. Bagnell. A reduction of imitation learning and structured prediction to no-regret online learning. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pages 627–635. JMLR Workshop and Conference Proceedings, 2011.
- E. Todorov, T. Erez, and Y. Tassa. Mujoco: A physics engine for model-based control. In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 5026–5033. IEEE, 2012.