

Linear Convergence of Natural Policy Gradient Methods with Log-Linear Policies

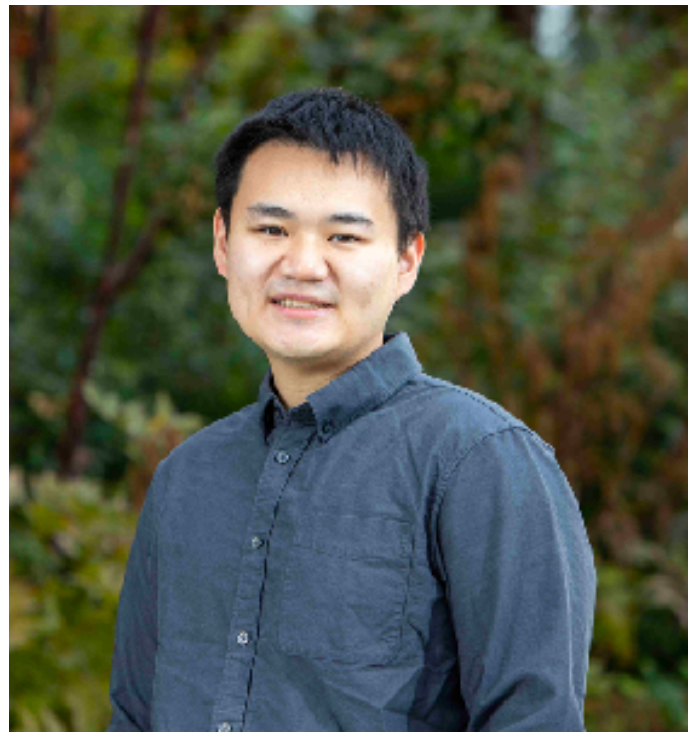
Rui Yuan^{1, 4}, Simon S. Du², Robert M. Gower³, Alessandro Lazaric¹, Lin Xiao¹

¹Meta AI, ²University of Washington, ³Flatiron Institute, ⁴Télécom Paris

International Conference on Learning Representations (ICLR), 2023



Thank you to



Simon S. Du²



Robert M. Gower³



Alessandro Lazaric¹



Lin Xiao¹

¹Meta AI ²University of Washington ³Flatiron Institute

Impressive Reinforcement Learning Results

Impressive Reinforcement Learning Results

Board Game

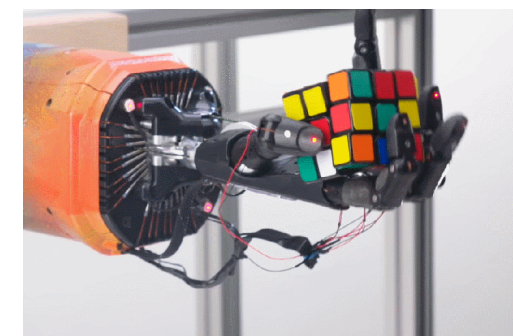
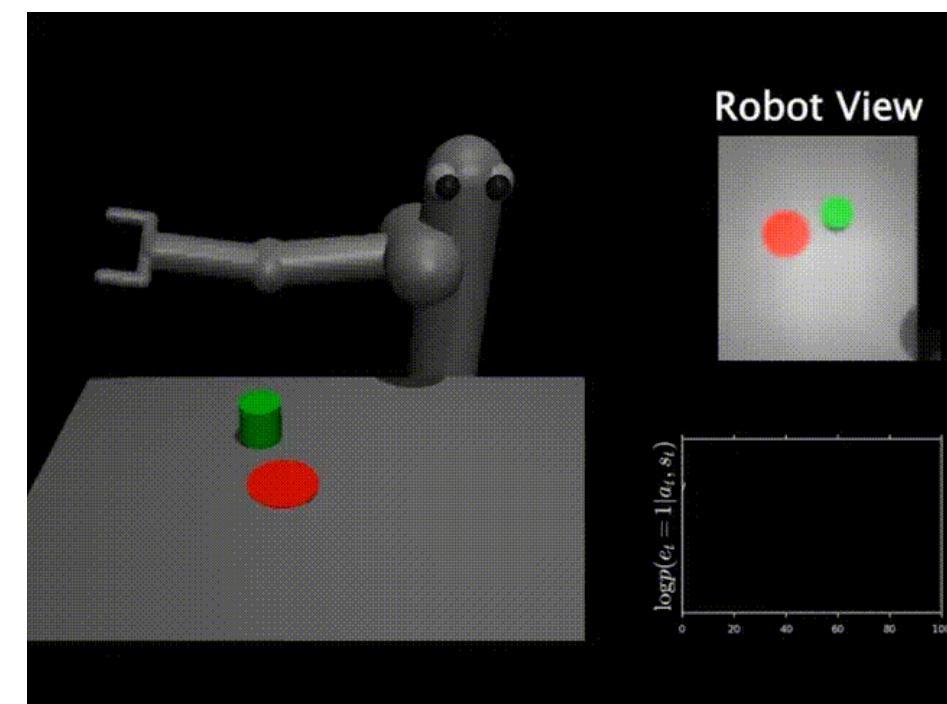


Impressive Reinforcement Learning Results

Board Game



Robotic Manipulation

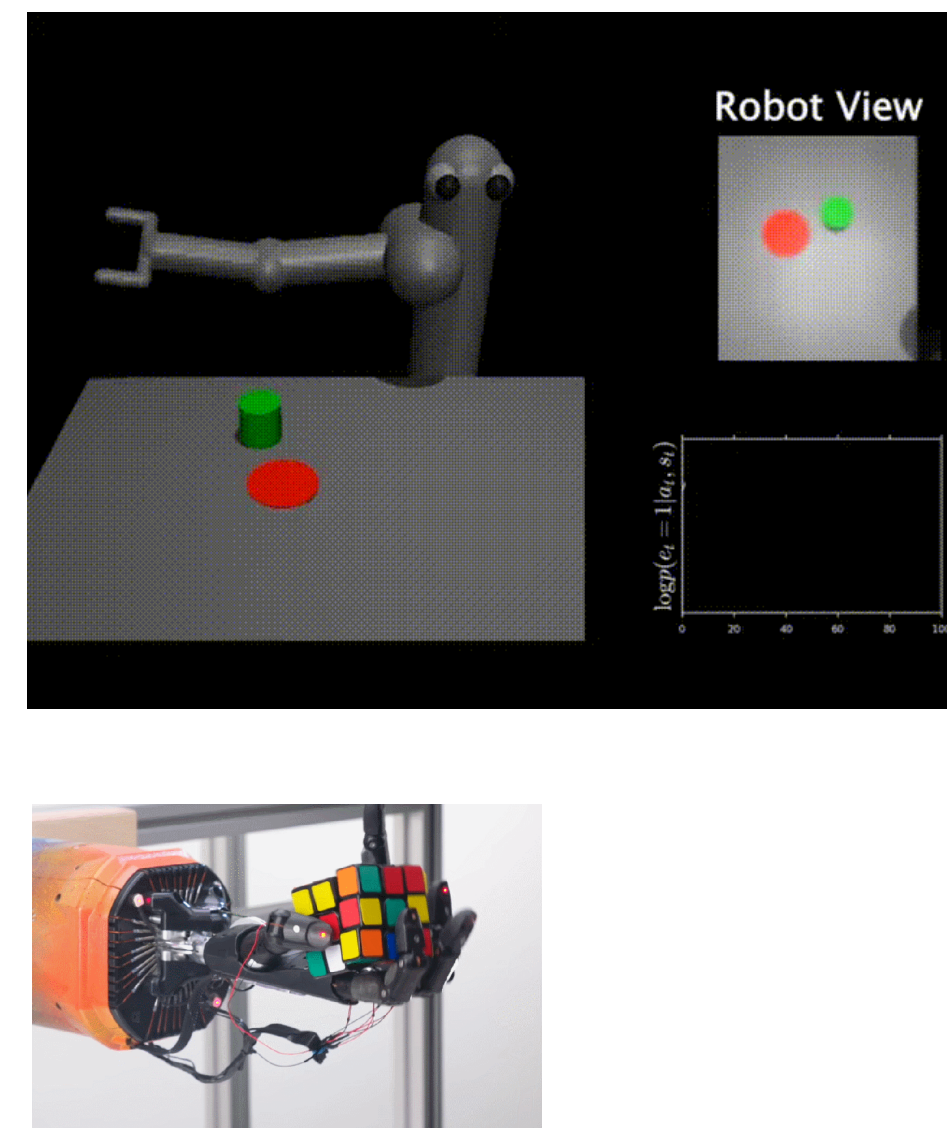


Impressive Reinforcement Learning Results

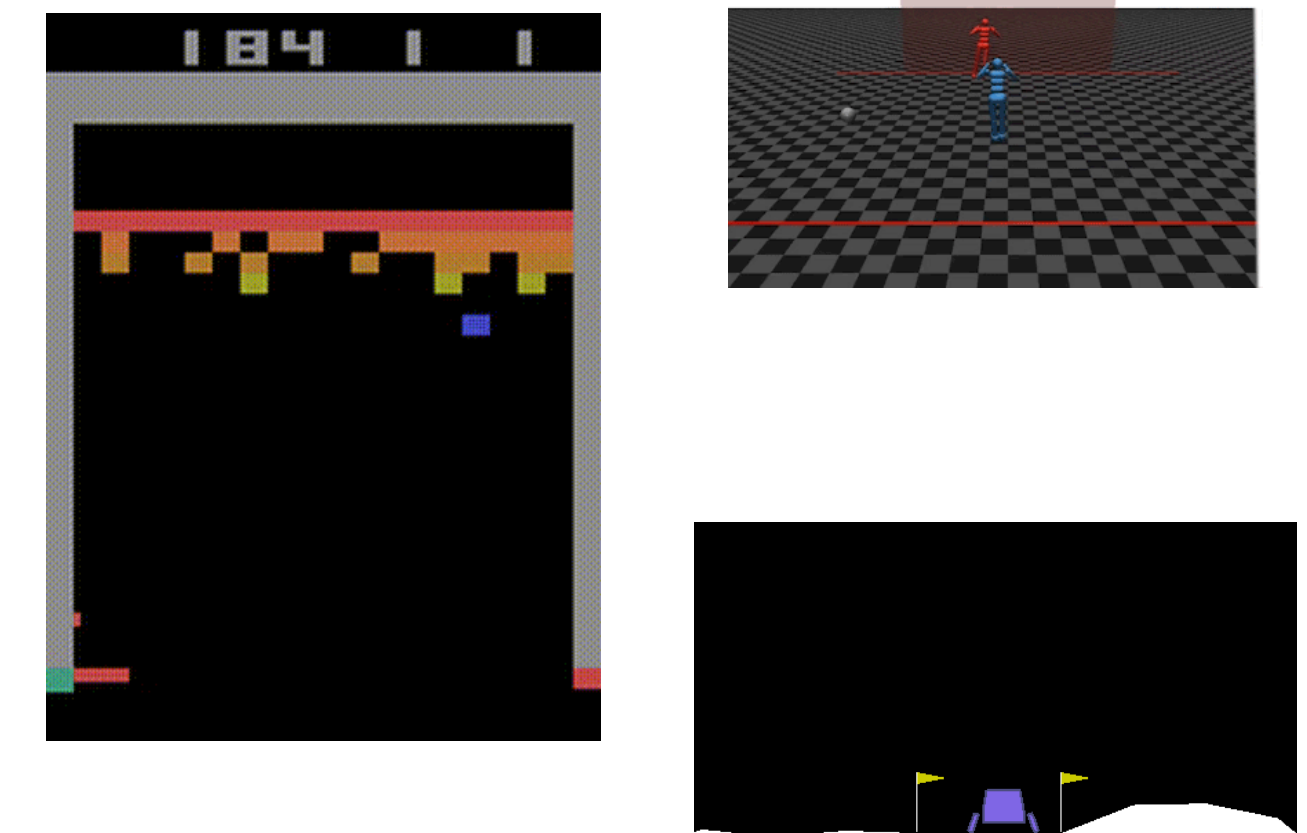
Board Game



Robotic Manipulation



Game Playing

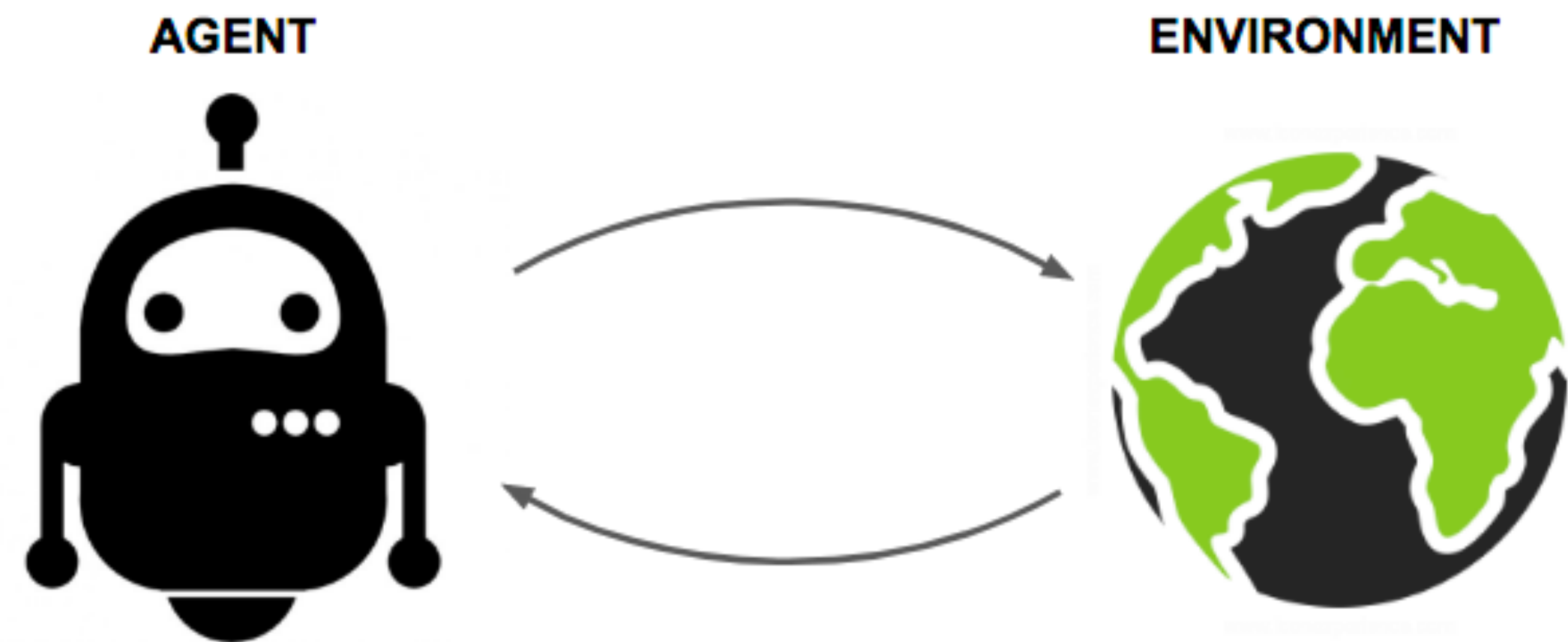


Reinforcement Learning (RL)

Sequential decision making problems

Reinforcement Learning (RL)

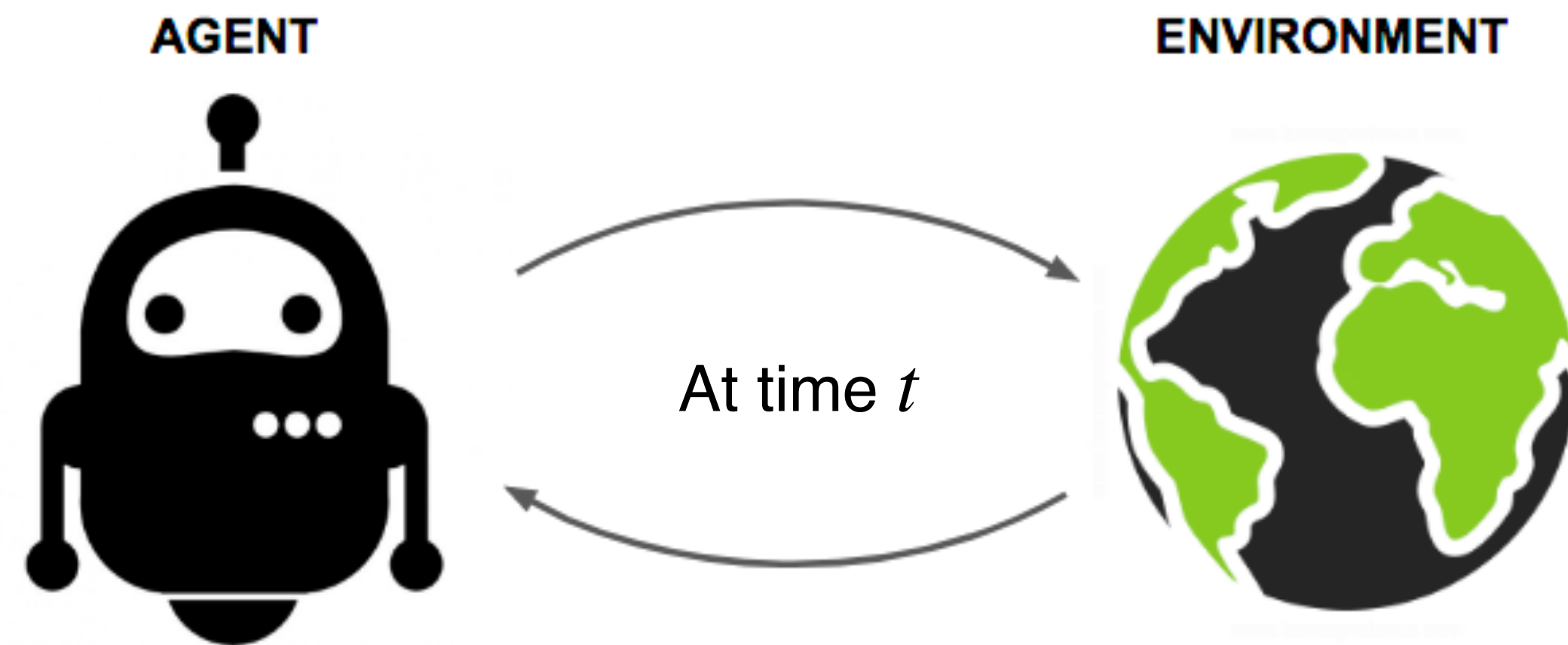
Sequential decision making problems



Markov decision Process (MDP)

Reinforcement Learning (RL)

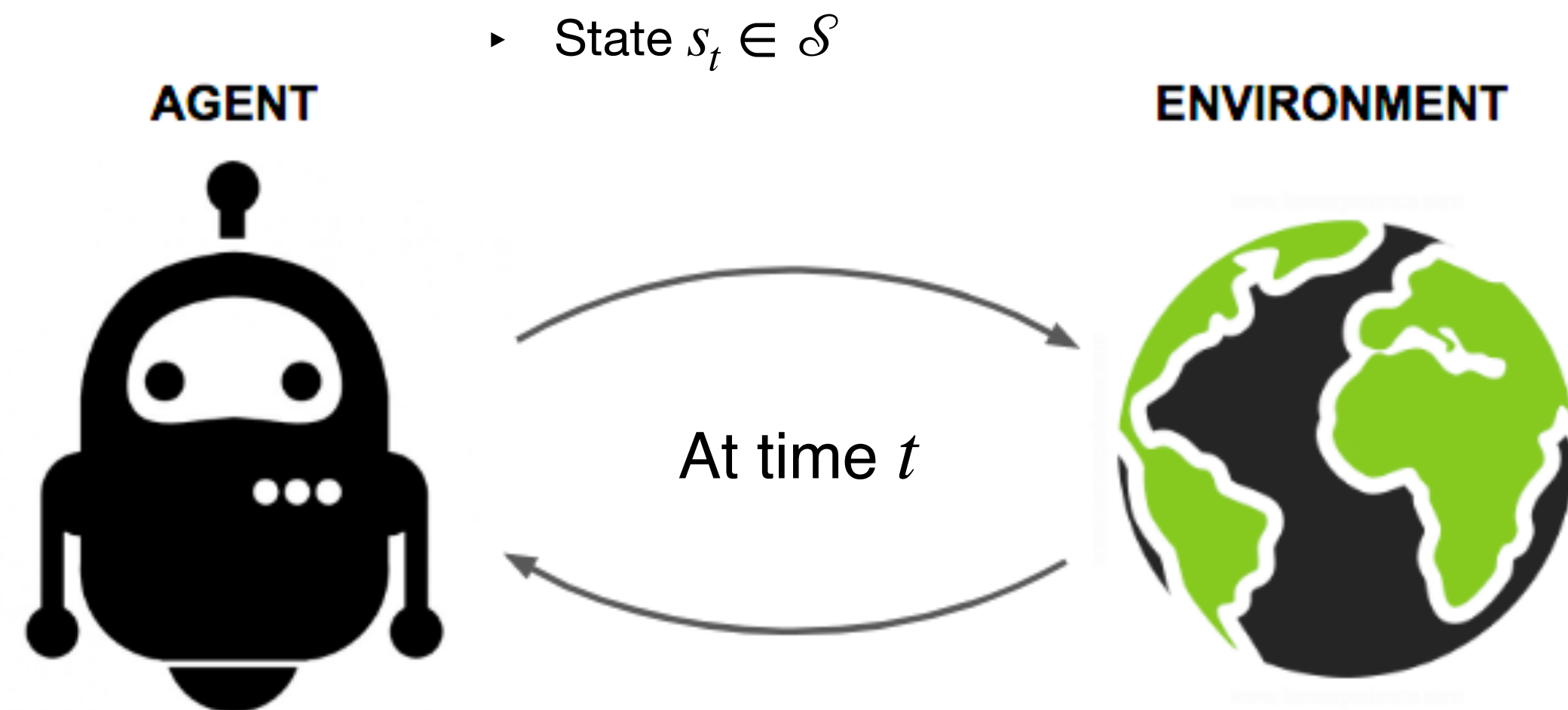
Sequential decision making problems



Markov decision Process (MDP)

Reinforcement Learning (RL)

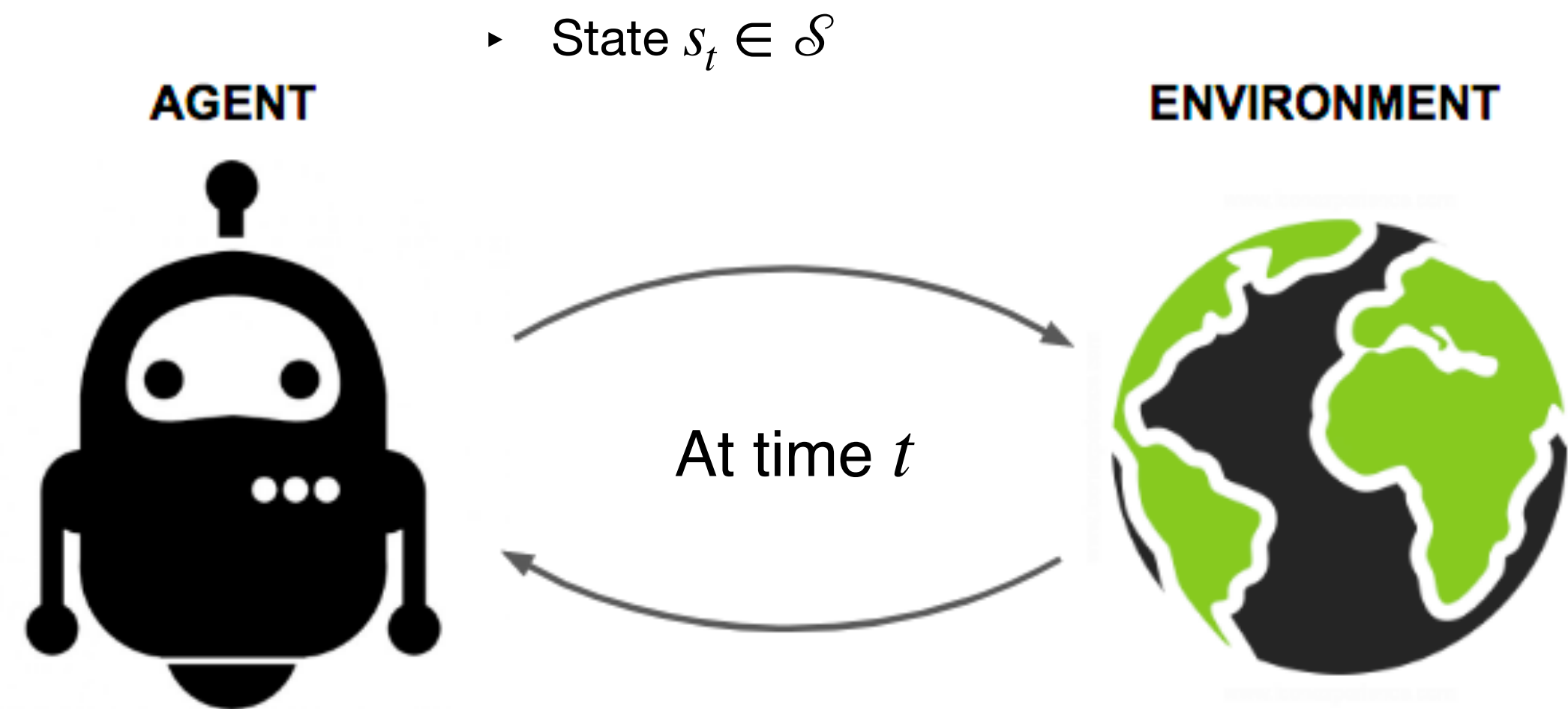
Sequential decision making problems



Markov decision Process (MDP)

Reinforcement Learning (RL)

Sequential decision making problems

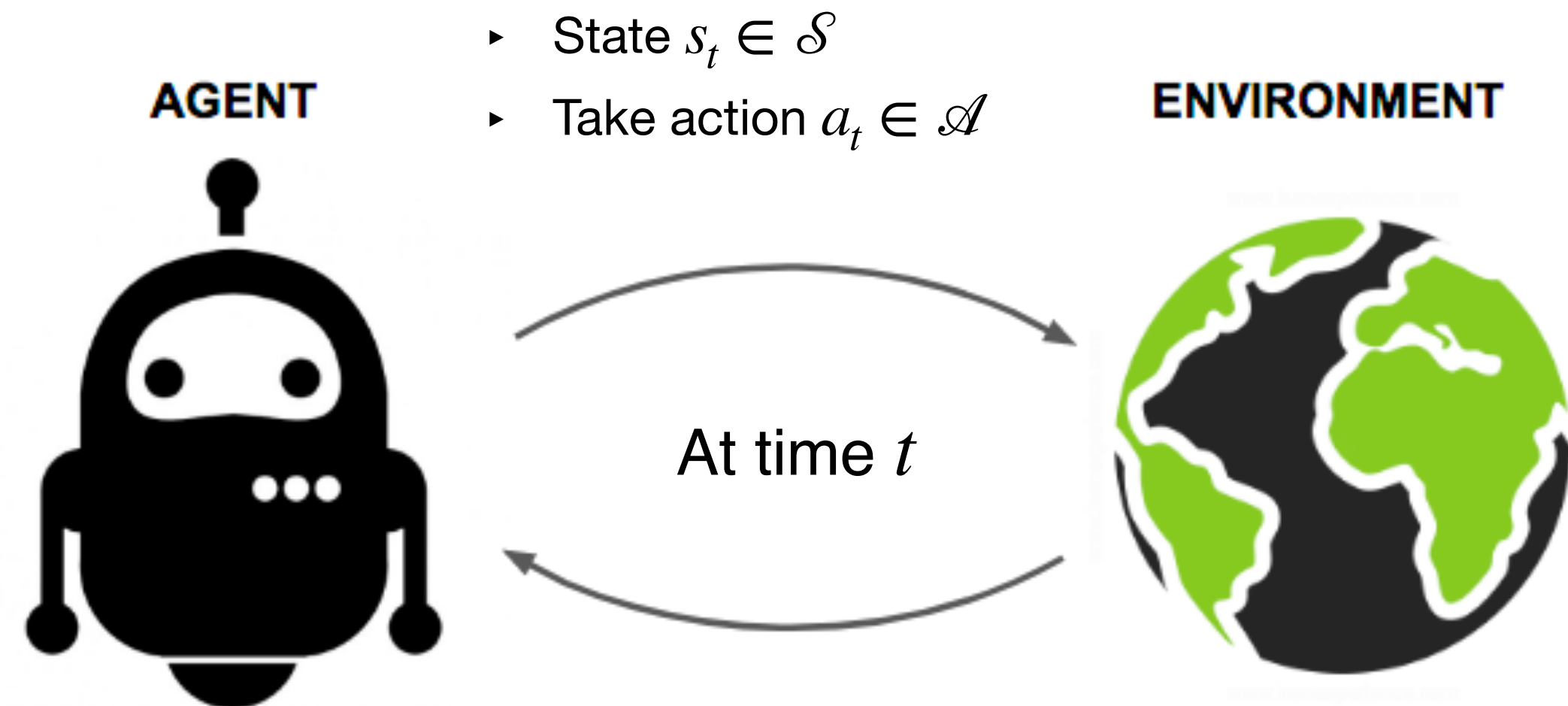


Markov decision Process (MDP)

- State space \mathcal{S}

Reinforcement Learning (RL)

Sequential decision making problems

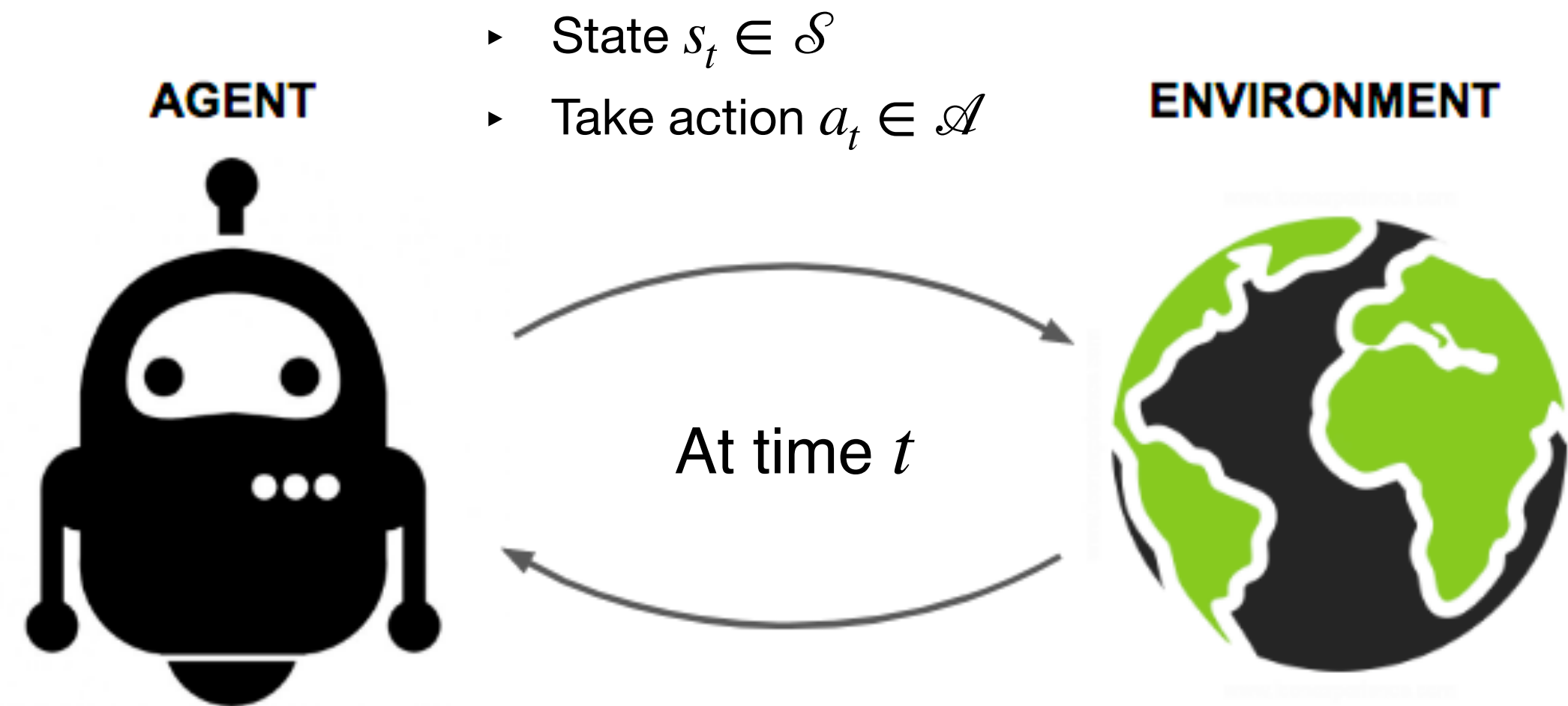


Markov decision Process (MDP)

- State space \mathcal{S}

Reinforcement Learning (RL)

Sequential decision making problems

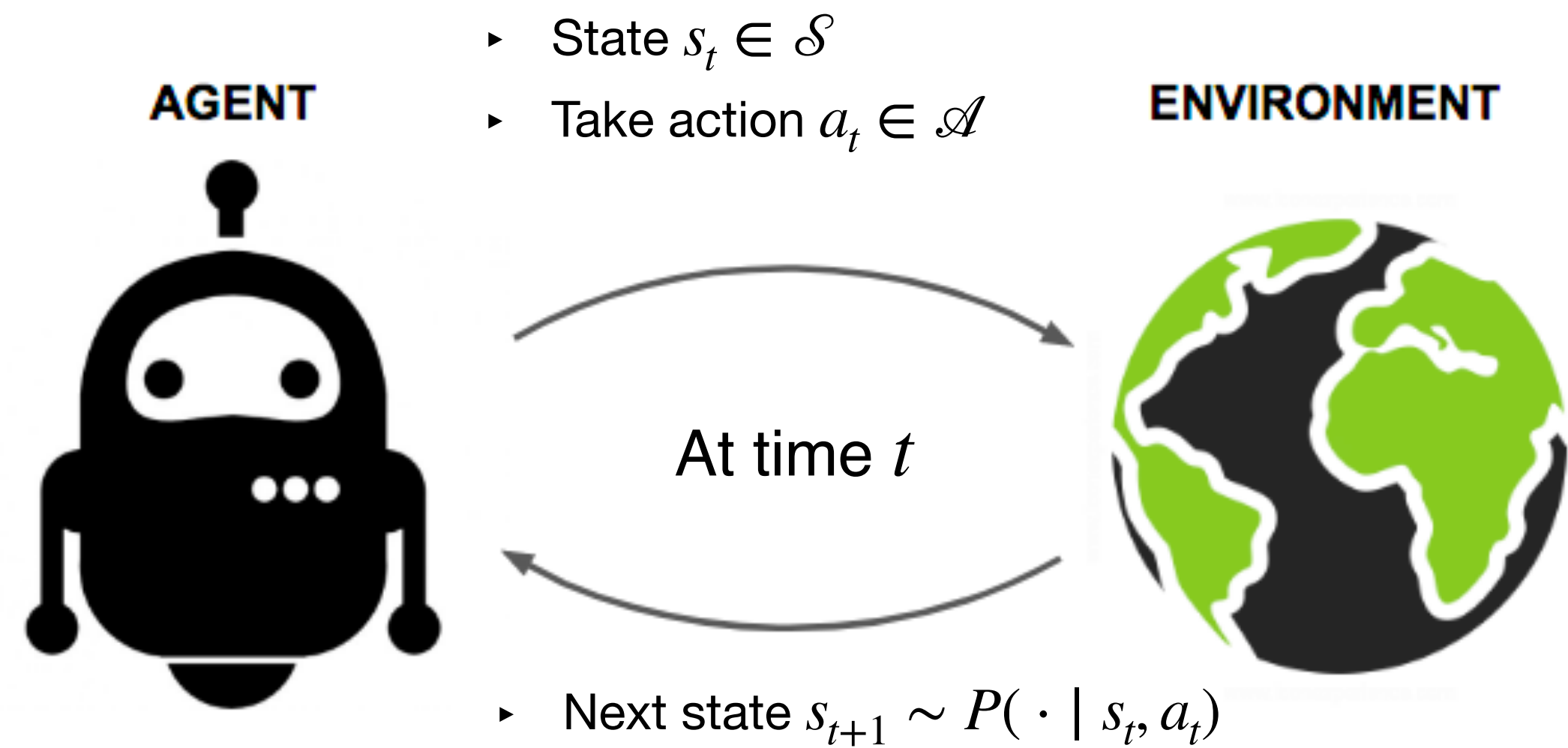


Markov decision Process (MDP)

- State space \mathcal{S}
- Action space \mathcal{A}

Reinforcement Learning (RL)

Sequential decision making problems

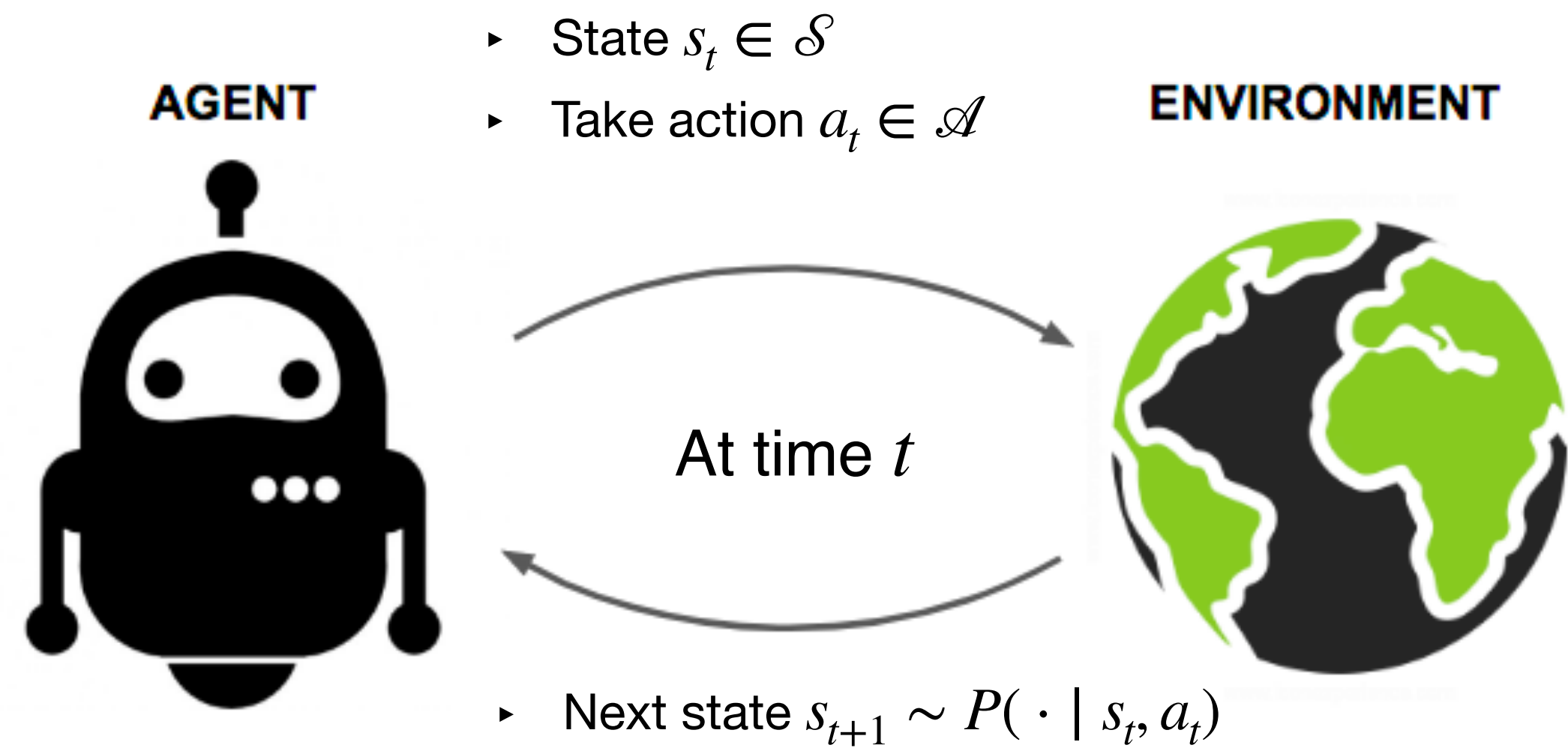


Markov decision Process (MDP)

- State space \mathcal{S}
- Action space \mathcal{A}

Reinforcement Learning (RL)

Sequential decision making problems

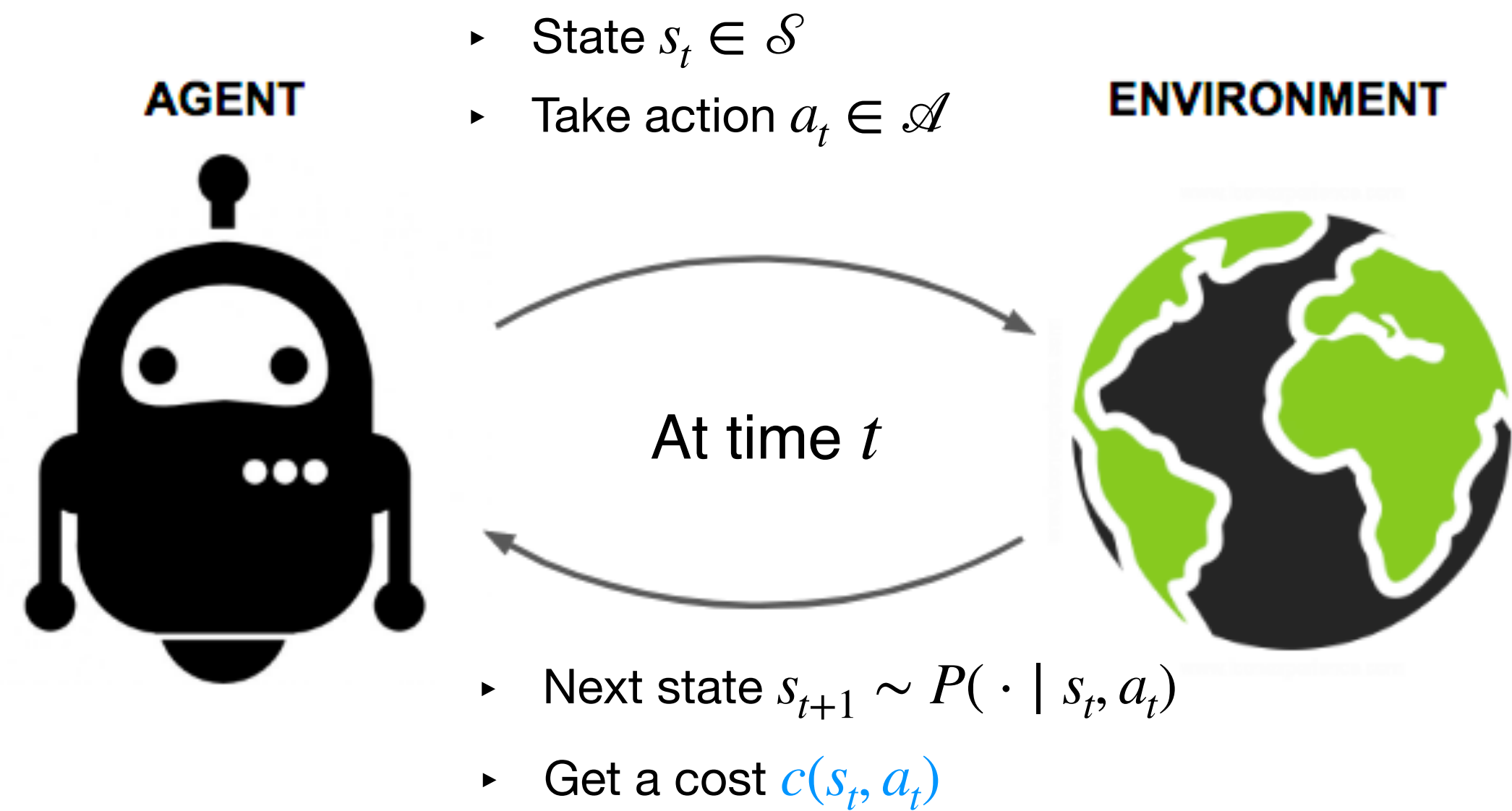


Markov decision Process (MDP)

- State space \mathcal{S}
- Action space \mathcal{A}
- Transition probabilities P

Reinforcement Learning (RL)

Sequential decision making problems

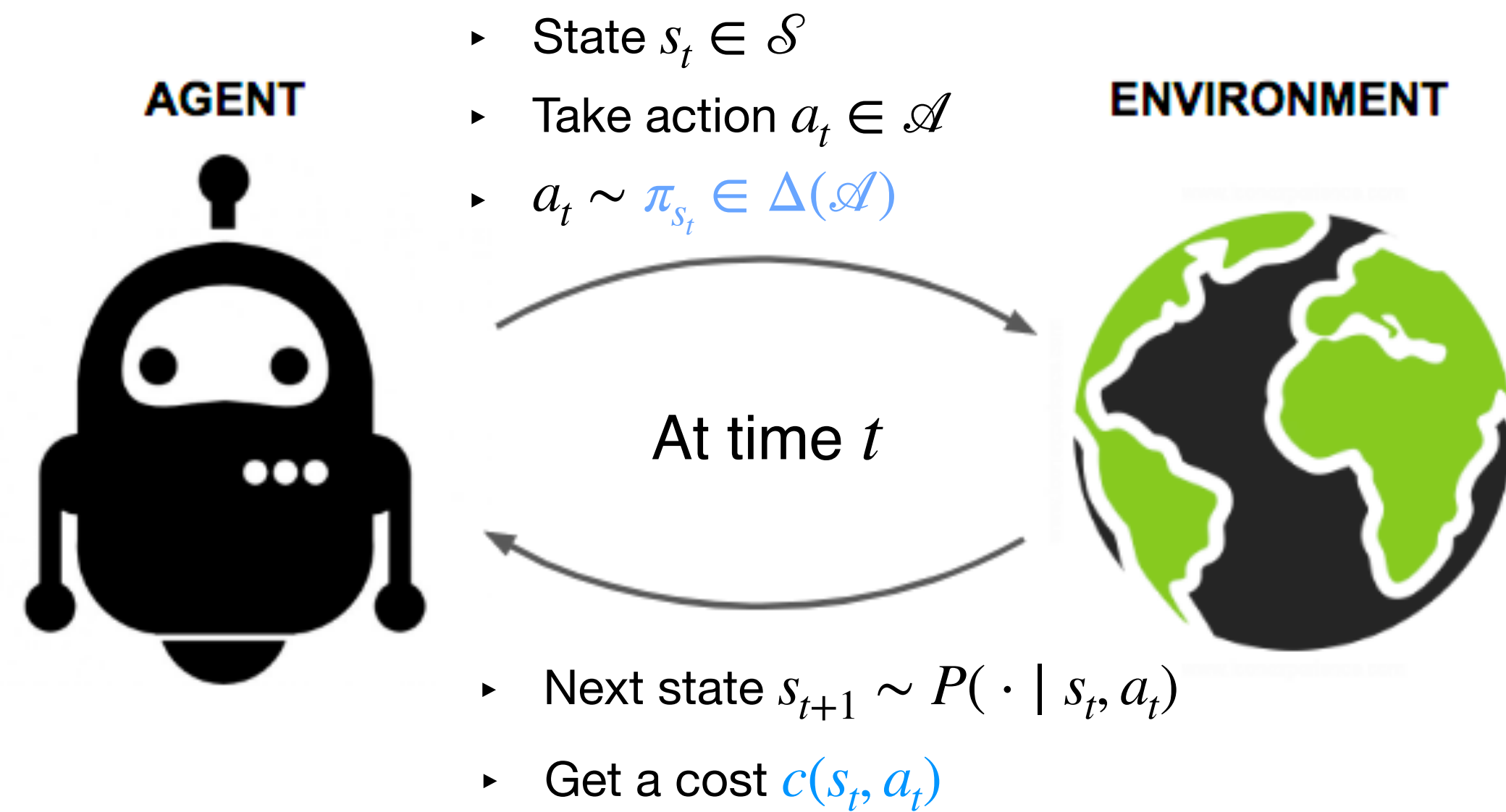


Markov decision Process (MDP)

- State space \mathcal{S}
- Action space \mathcal{A}
- Transition probabilities P

Reinforcement Learning (RL)

Sequential decision making problems



Markov decision Process (MDP)

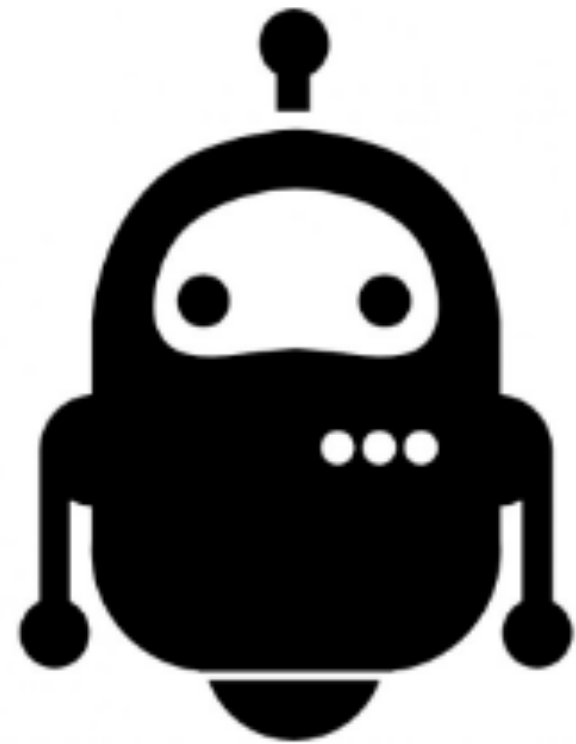
- State space \mathcal{S}
- Action space \mathcal{A}
- Transition probabilities P

Reinforcement Learning (RL)

Sequential decision making problems

Policy $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$,

AGENT



- ▶ State $s_t \in \mathcal{S}$
- ▶ Take action $a_t \in \mathcal{A}$
- ▶ $a_t \sim \pi_{s_t} \in \Delta(\mathcal{A})$

ENVIRONMENT



At time t

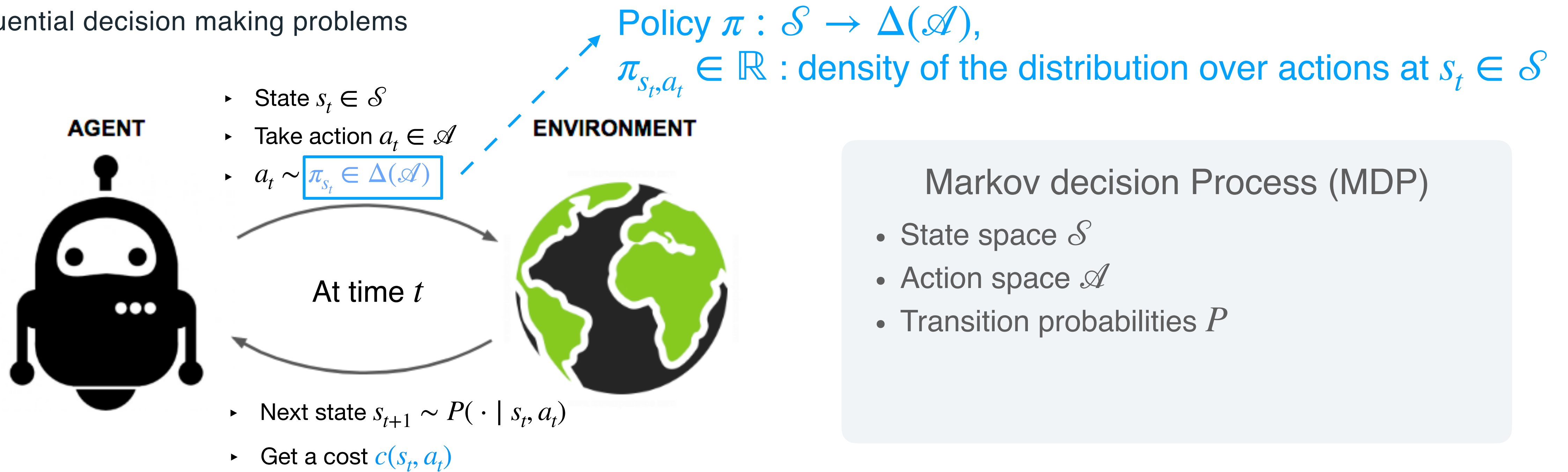
- ▶ Next state $s_{t+1} \sim P(\cdot | s_t, a_t)$
- ▶ Get a cost $c(s_t, a_t)$

Markov decision Process (MDP)

- State space \mathcal{S}
- Action space \mathcal{A}
- Transition probabilities P

Reinforcement Learning (RL)

Sequential decision making problems

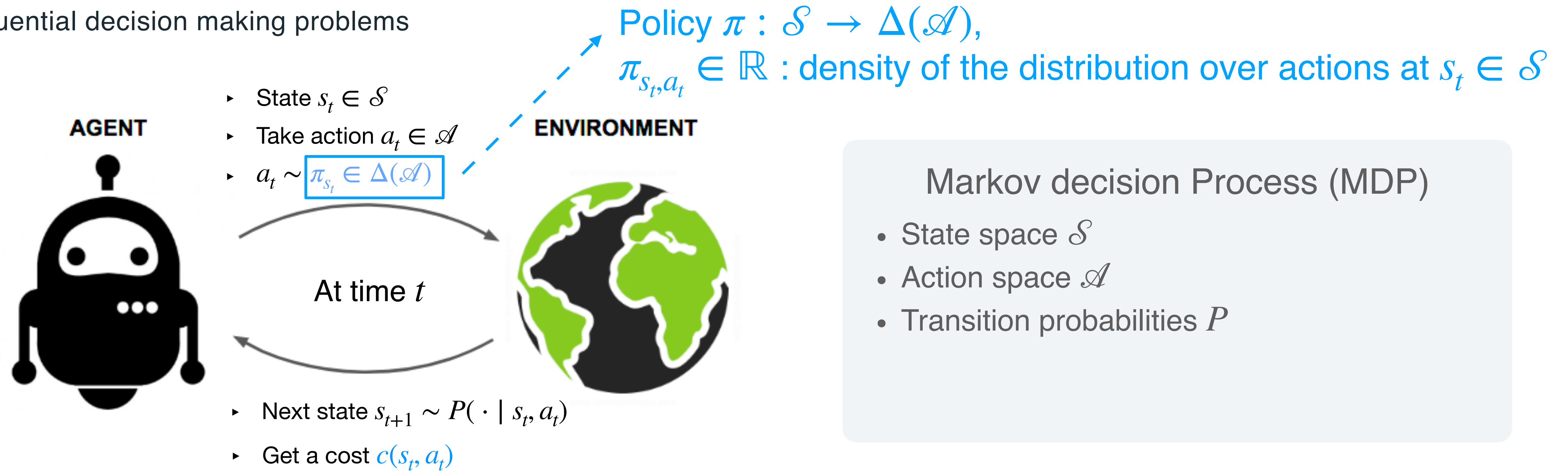


Markov decision Process (MDP)

- State space \mathcal{S}
- Action space \mathcal{A}
- Transition probabilities P

Reinforcement Learning (RL)

Sequential decision making problems

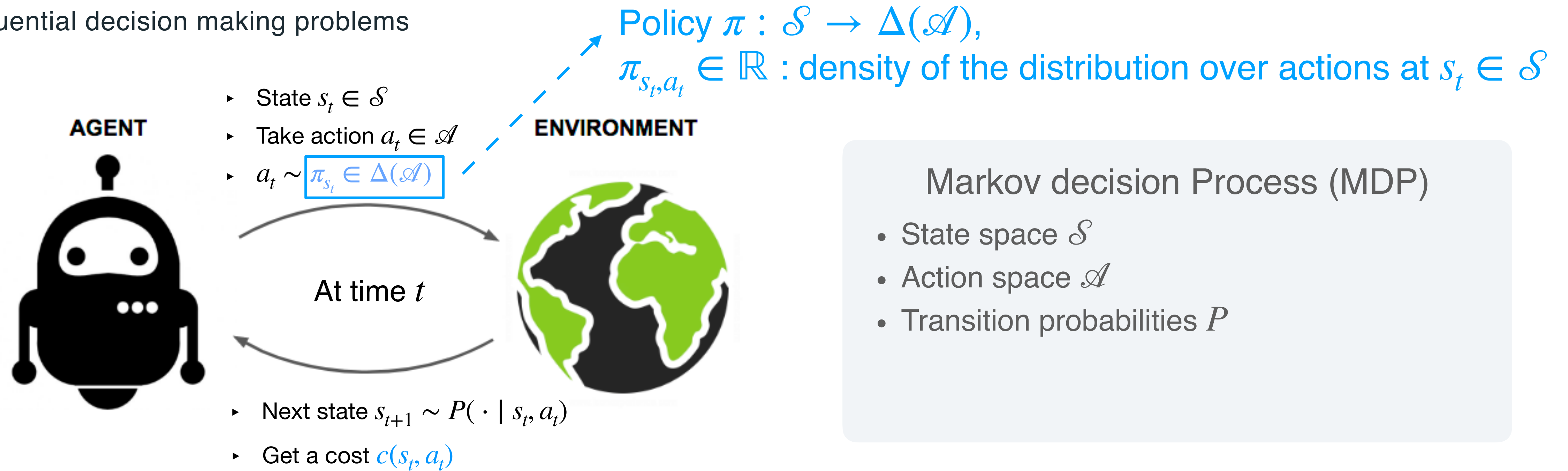


Solve an MDP to minimize total expected cost (a.k.a. policy optimization)

$$\arg \min_{\pi} V_{\rho}(\pi) := \mathbb{E}_{s_0 \sim \rho, a_t \sim \pi_{s_t}, s_{t+1} \sim P(\cdot | s_t, a_t)} \left[\sum_{t=0}^{\infty} \gamma^t c(s_t, a_t) \right]$$

Reinforcement Learning (RL)

Sequential decision making problems

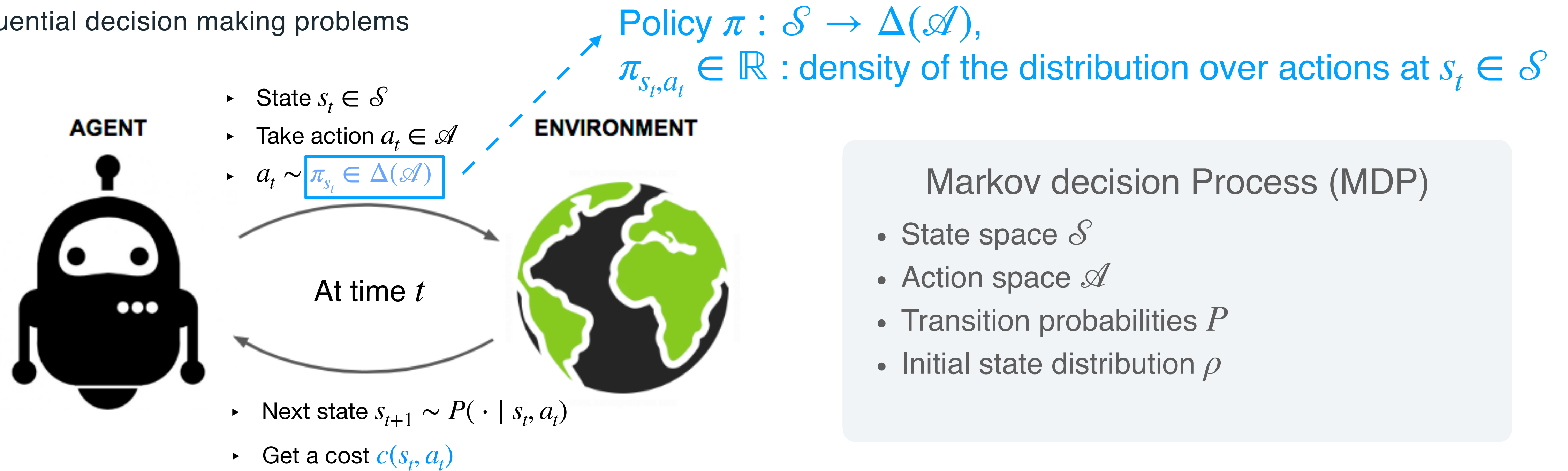


Solve an MDP to minimize total expected cost (a.k.a. policy optimization)

$$\arg \min_{\pi} V_{\rho}(\pi) := \mathbb{E}_{s_0 \sim \rho, a_t \sim \pi_{s_t}, s_{t+1} \sim P(\cdot | s_t, a_t)} \left[\sum_{t=0}^{\infty} \gamma^t c(s_t, a_t) \right] \rightarrow \text{Cost function}$$

Reinforcement Learning (RL)

Sequential decision making problems

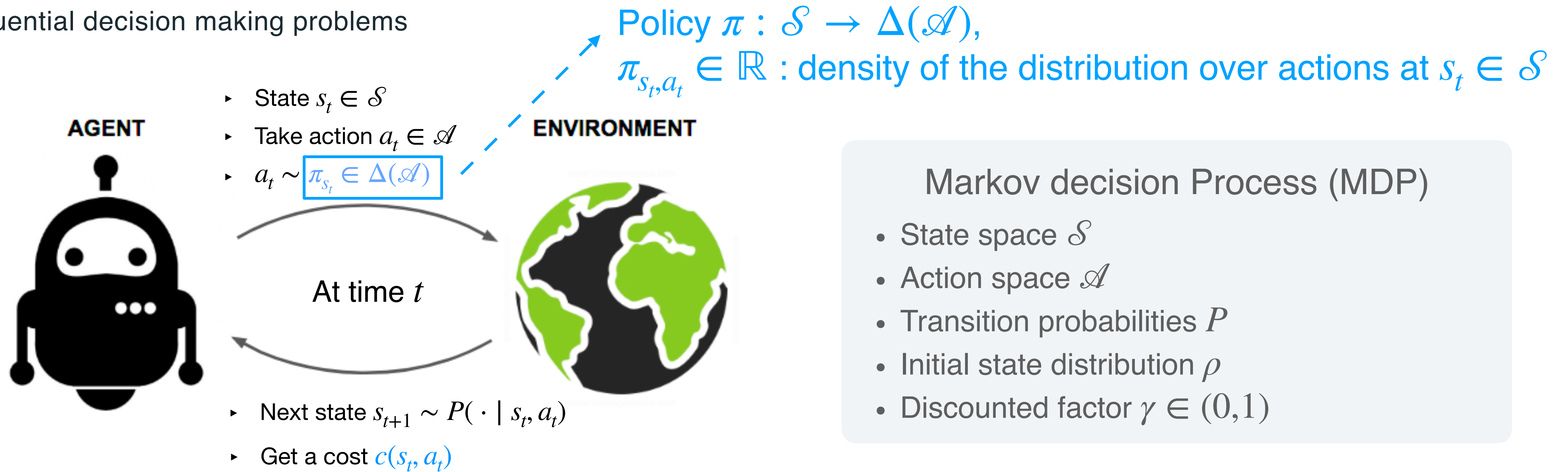


Solve an MDP to **minimize total expected cost** (a.k.a. **policy optimization**)

$$\arg \min_{\pi} V_{\rho}(\pi) := \mathbb{E}_{s_0 \sim \rho, a_t \sim \pi_{s_t}, s_{t+1} \sim P(\cdot | s_t, a_t)} \left[\sum_{t=0}^{\infty} \gamma^t c(s_t, a_t) \right] \rightarrow \text{Cost function}$$

Reinforcement Learning (RL)

Sequential decision making problems



Markov decision Process (MDP)

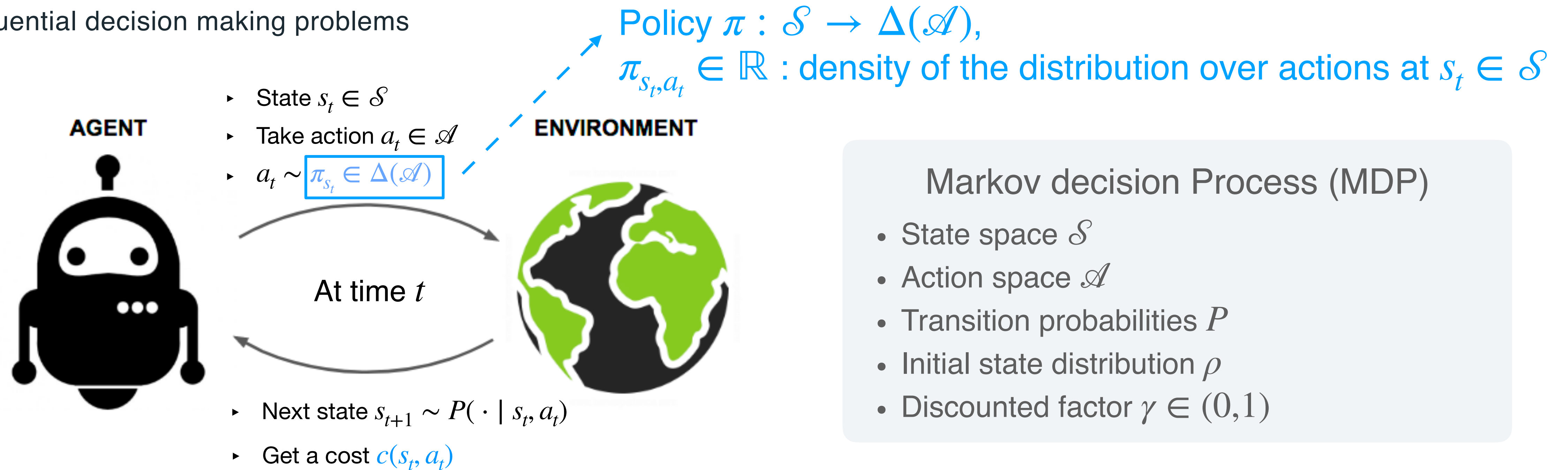
- State space \mathcal{S}
- Action space \mathcal{A}
- Transition probabilities P
- Initial state distribution ρ
- Discounted factor $\gamma \in (0, 1)$

Solve an MDP to minimize total expected cost (a.k.a. policy optimization)

$$\arg \min_{\pi} V_{\rho}(\pi) := \mathbb{E}_{s_0 \sim \rho, a_t \sim \pi_{s_t}, s_{t+1} \sim P(\cdot | s_t, a_t)} \left[\sum_{t=0}^{\infty} \gamma^t c(s_t, a_t) \right] \rightarrow \text{Cost function}$$

Reinforcement Learning (RL)

Sequential decision making problems

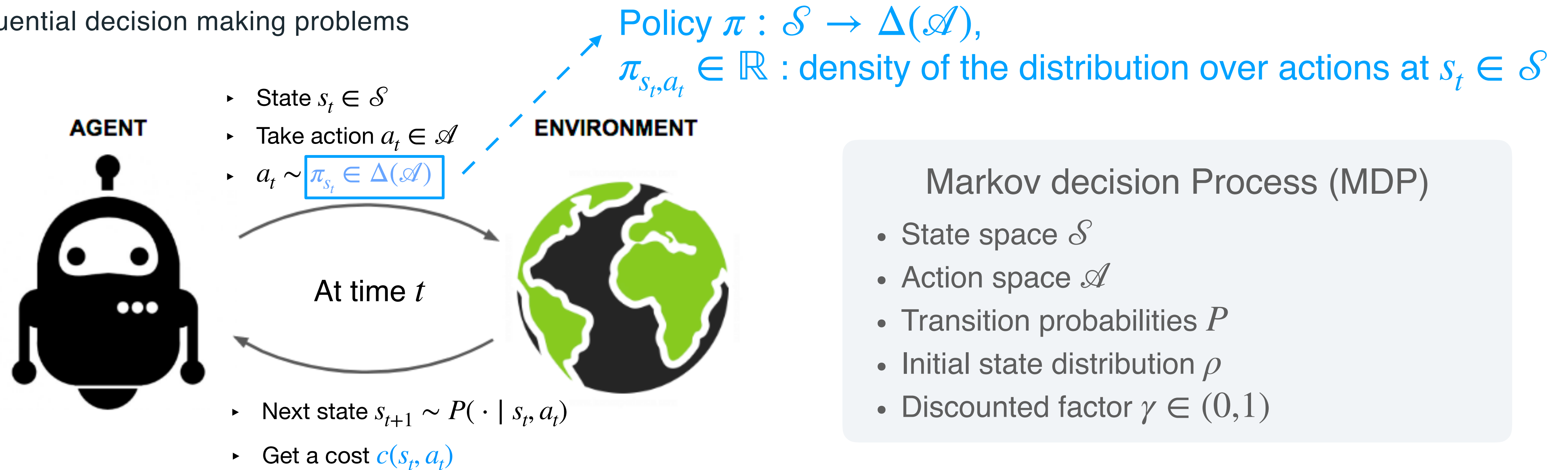


Solve an MDP to minimize total expected cost (a.k.a. policy optimization)

$$\arg \min_{\theta \in \mathbb{R}^d} V_{\rho}(\theta) := \mathbb{E}_{s_0 \sim \rho, a_t \sim \pi_{s_t}(\theta), s_{t+1} \sim P(\cdot | s_t, a_t)} \left[\sum_{t=0}^{\infty} \gamma^t c(s_t, a_t) \right]$$

Reinforcement Learning (RL)

Sequential decision making problems



Solve an MDP to minimize total expected cost (a.k.a. policy optimization)

$$\arg \min_{\theta \in \mathbb{R}^d} V_{\rho}(\theta) := \mathbb{E}_{s_0 \sim \rho, a_t \sim \pi_{s_t}(\theta), s_{t+1} \sim P(\cdot | s_t, a_t)} \left[\sum_{t=0}^{\infty} \gamma^t c(s_t, a_t) \right]$$

Objective: $\arg \min_{\theta \in \mathbb{R}^d} V_{\rho}(\theta)$

Objective: $\arg \min_{\theta \in \mathbb{R}^d} V_\rho(\theta)$

$$\theta^{(k+1)} = \theta^{(k)} - \eta_k \nabla_{\theta} V_\rho(\theta^{(k)})$$

Objective: $\arg \min_{\theta \in \mathbb{R}^d} V_\rho(\theta)$

$$\theta^{(k+1)} = \theta^{(k)} - \eta_k \nabla_{\theta} V_\rho(\theta^{(k)})$$

Step size

Objective: $\arg \min_{\theta \in \mathbb{R}^d} V_\rho(\theta)$

$$\theta^{(k+1)} = \theta^{(k)} - \eta_k \nabla_{\theta} V_\rho(\theta^{(k)})$$

Step size

Gradient of $V_\rho(\theta)$

Policy gradient (PG) methods

Objective: $\arg \min_{\theta \in \mathbb{R}^d} V_{\rho}(\theta)$

$$\theta^{(k+1)} = \theta^{(k)} - \eta_k \nabla_{\theta} V_{\rho}(\theta^{(k)})$$

Step size

Gradient of $V_{\rho}(\theta)$

Policy gradient (PG) methods

Objective: $\arg \min_{\theta \in \mathbb{R}^d} V_{\rho}(\theta)$

- Simplicity

$$\theta^{(k+1)} = \theta^{(k)} - \eta_k \nabla_{\theta} V_{\rho}(\theta^{(k)})$$

Step size

Gradient of $V_{\rho}(\theta)$

The diagram shows the update equation $\theta^{(k+1)} = \theta^{(k)} - \eta_k \nabla_{\theta} V_{\rho}(\theta^{(k)})$ enclosed in a black rectangular box. A blue arrow points from the text 'Step size' below to the η_k term. Another blue arrow points from the text 'Gradient of $V_{\rho}(\theta)$ ' to the $\nabla_{\theta} V_{\rho}(\theta^{(k)})$ term.

Policy gradient (PG) methods

Objective: $\arg \min_{\theta \in \mathbb{R}^d} V_{\rho}(\theta)$

- Simplicity
 - Easy to implement and use in practice

$$\theta^{(k+1)} = \theta^{(k)} - \eta_k \nabla_{\theta} V_{\rho}(\theta^{(k)})$$

Step size

Gradient of $V_{\rho}(\theta)$

Policy gradient (PG) methods

Objective: $\arg \min_{\theta \in \mathbb{R}^d} V_{\rho}(\theta)$

$$\theta^{(k+1)} = \theta^{(k)} - \eta_k \nabla_{\theta} V_{\rho}(\theta^{(k)})$$

Step size

Gradient of $V_{\rho}(\theta)$

- Simplicity
 - Easy to implement and use in practice
 - Can solve a wide range of problems (e.g. partially-observable environments)

Policy gradient (PG) methods

Objective: $\arg \min_{\theta \in \mathbb{R}^d} V_{\rho}(\theta)$

$$\theta^{(k+1)} = \theta^{(k)} - \eta_k \nabla_{\theta} V_{\rho}(\theta^{(k)})$$

Step size

Gradient of $V_{\rho}(\theta)$

- Simplicity
 - Easy to implement and use in practice
 - Can solve a wide range of problems (e.g. partially-observable environments)
- Versatility

Policy gradient (PG) methods

Objective: $\arg \min_{\theta \in \mathbb{R}^d} V_{\rho}(\theta)$

$$\theta^{(k+1)} = \theta^{(k)} - \eta_k \nabla_{\theta} V_{\rho}(\theta^{(k)})$$

Step size

Gradient of $V_{\rho}(\theta)$

- Simplicity

- Easy to implement and use in practice

- Can solve a wide range of problems (e.g. partially-observable environments)

- Versatility

- REINFORCE [Williams, 1992], PGT, GPOMDP [Baxter and Bartlett, 2001], actor-critic [Konda and Tsitsiklis, 2000]

Policy gradient (PG) methods

Objective: $\arg \min_{\theta \in \mathbb{R}^d} V_{\rho}(\theta)$

$$\theta^{(k+1)} = \theta^{(k)} - \eta_k \nabla_{\theta} V_{\rho}(\theta^{(k)})$$

Step size

Gradient of $V_{\rho}(\theta)$

- Simplicity

- Easy to implement and use in practice
- Can solve a wide range of problems (e.g. partially-observable environments)

- Versatility

- REINFORCE [Williams, 1992], PGT, GPOMDP [Baxter and Bartlett, 2001], actor-critic [Konda and Tsitsiklis, 2000]
- Natural PG [Kakade, 2001], policy mirror descent [Lan, 2022; Xiao, 2022], variance reduction techniques [Papini et al., 2018; Shen et al., 2019; Xu et al., 2020; Huang et al., 2020]

Policy gradient (PG) methods

Objective: $\arg \min_{\theta \in \mathbb{R}^d} V_{\rho}(\theta)$

$$\theta^{(k+1)} = \theta^{(k)} - \eta_k \nabla_{\theta} V_{\rho}(\theta^{(k)})$$

Annotations:
• η_k : Step size
• $\nabla_{\theta} V_{\rho}(\theta^{(k)})$: Gradient of $V_{\rho}(\theta)$

- Simplicity

- Easy to implement and use in practice
- Can solve a wide range of problems (e.g. partially-observable environments)

- Versatility

- REINFORCE [Williams, 1992], PGT, GPOMDP [Baxter and Bartlett, 2001], actor-critic [Konda and Tsitsiklis, 2000]
- Natural PG [Kakade, 2001], policy mirror descent [Lan, 2022; Xiao, 2022], variance reduction techniques [Papini et al., 2018; Shen et al., 2019; Xu et al., 2020; Huang et al., 2020]
- Trust-region (e.g. TRPO [Schulman et al., 2015;]), proximal (e.g. PPO [Schulman et al., 2017])

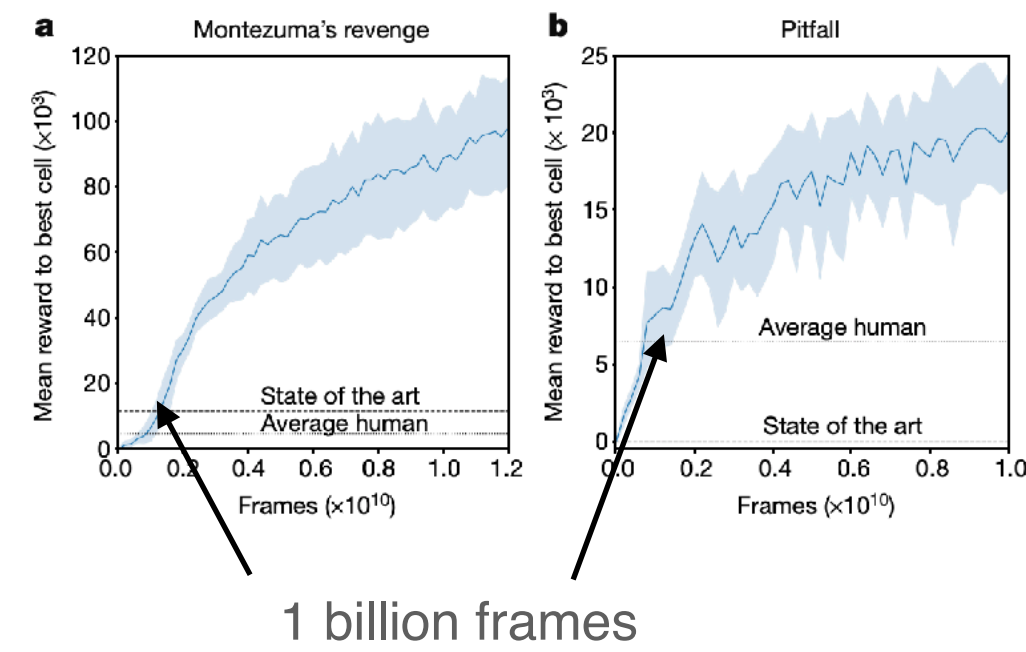
Context

Objective: $\arg \min_{\theta \in \mathbb{R}^d} V_{\rho}(\theta)$

Context

Objective: $\arg \min_{\theta \in \mathbb{R}^d} V_{\rho}(\theta)$

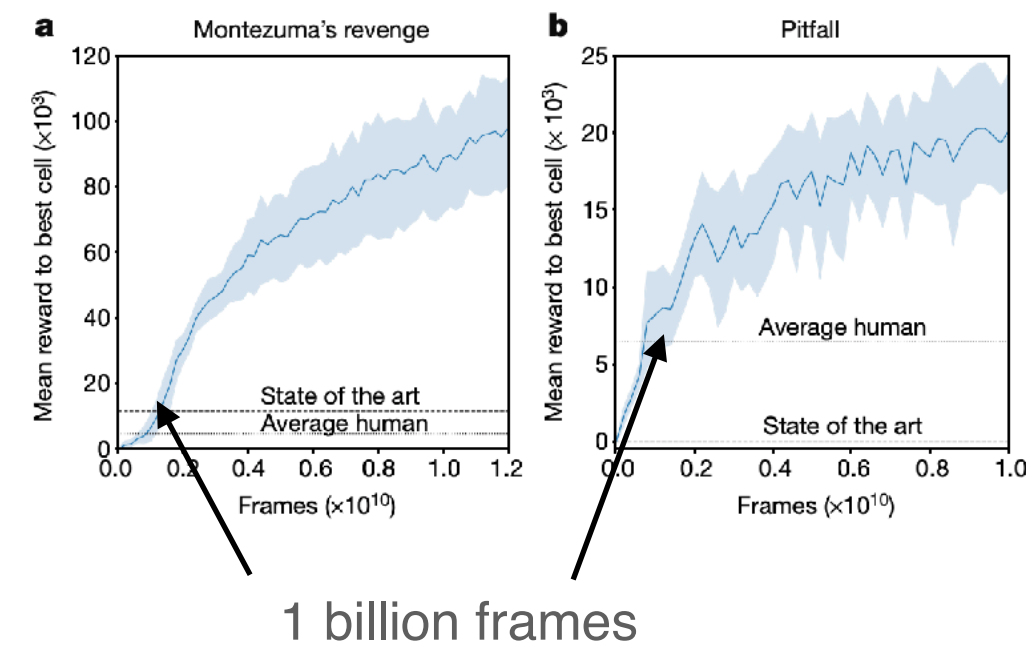
- Original PG is not **sample efficient**



Context

Objective: $\arg \min_{\theta \in \mathbb{R}^d} V_{\rho}(\theta)$

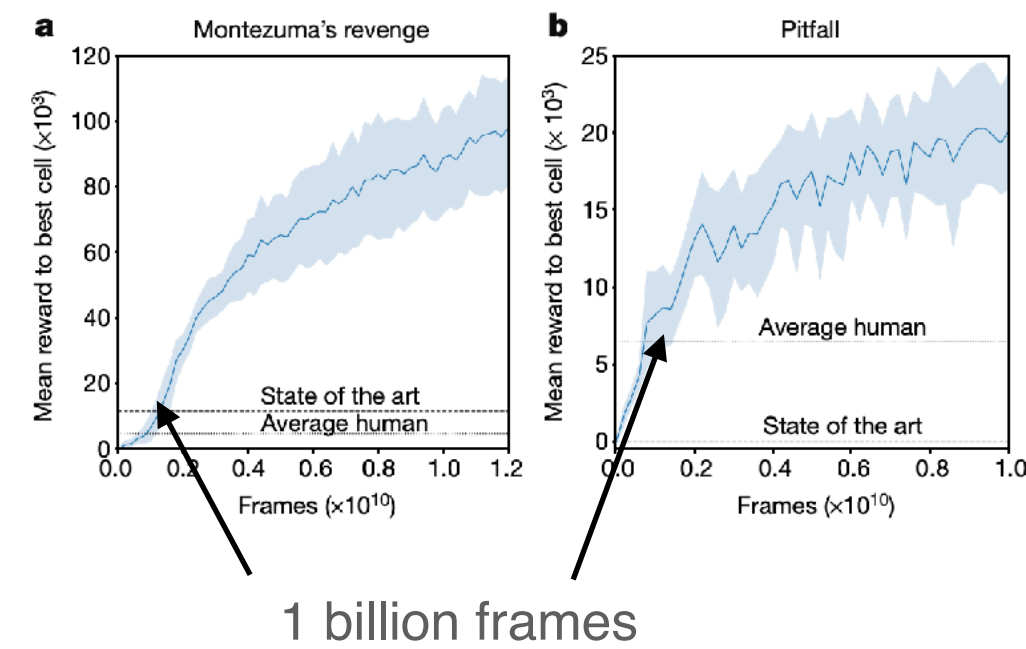
- Original PG is not **sample efficient**
- Natural PG (NPG)[Kakade, 2001] uses a **preconditioner** to improve PG direction



Context

Objective: $\arg \min_{\theta \in \mathbb{R}^d} V_{\rho}(\theta)$

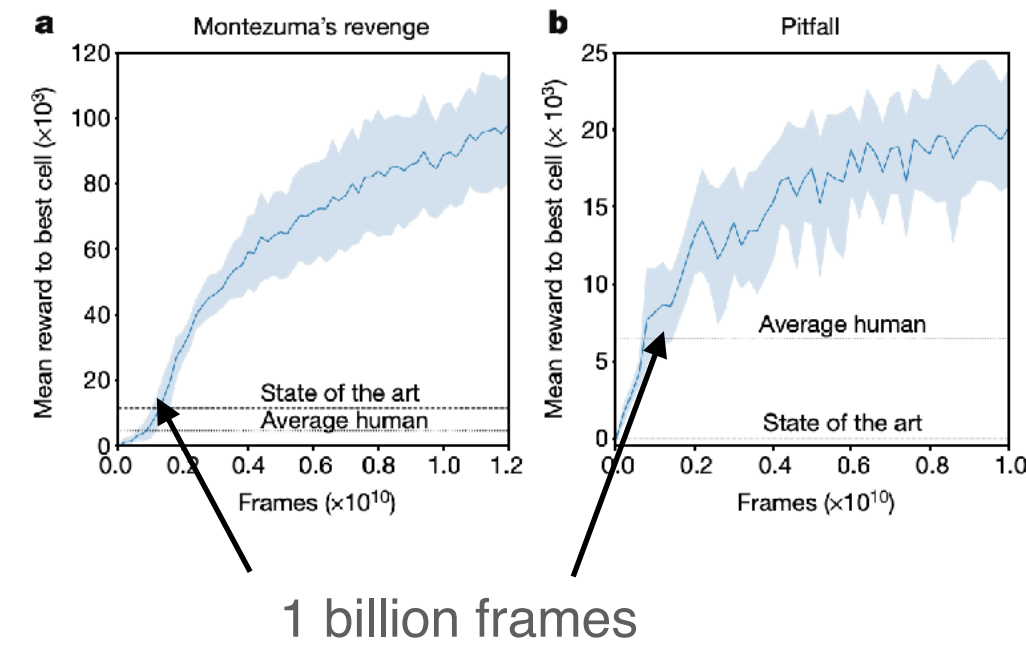
- Original PG is not **sample efficient**
- Natural PG (NPG)[Kakade, 2001] uses a **preconditioner** to improve PG direction
- NPG is the **building block** of several **state-of-the-art** algorithms (**TRPO, PPO**)



Context

Objective: $\arg \min_{\theta \in \mathbb{R}^d} V_{\rho}(\theta)$

- Original PG is not **sample efficient**

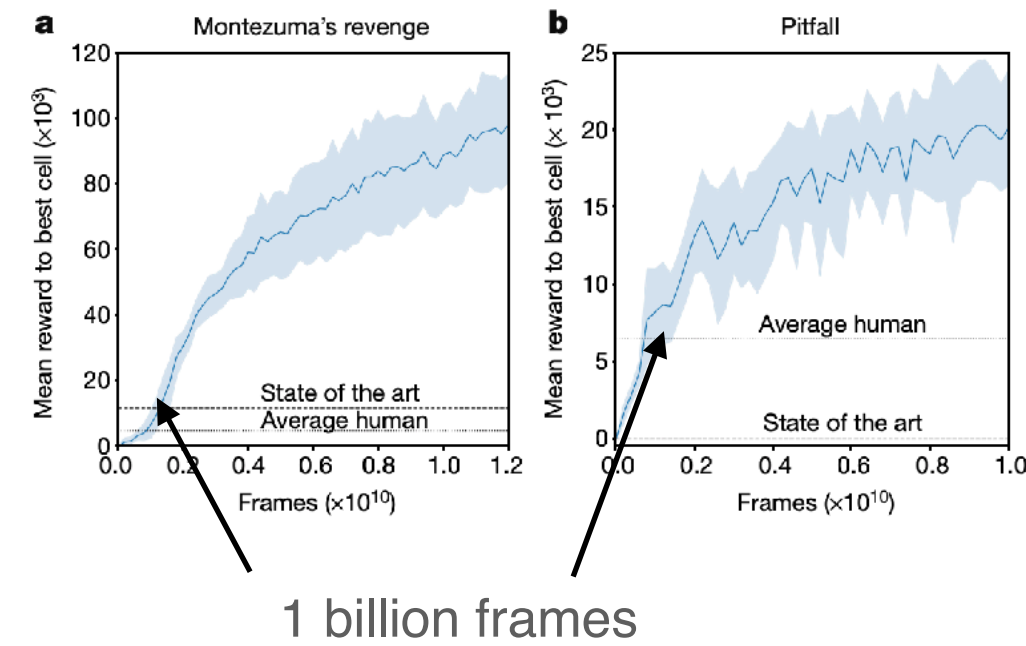


- Natural PG (NPG)[Kakade, 2001] uses a **preconditioner** to improve PG direction
- NPG is the **building block** of several **state-of-the-art** algorithms (**TRPO, PPO**)
- **Linear convergence** of NPG is established for tabular case [Xiao, 2022]

Context

Objective: $\arg \min_{\theta \in \mathbb{R}^d} V_{\rho}(\theta)$

- Original PG is not **sample efficient**

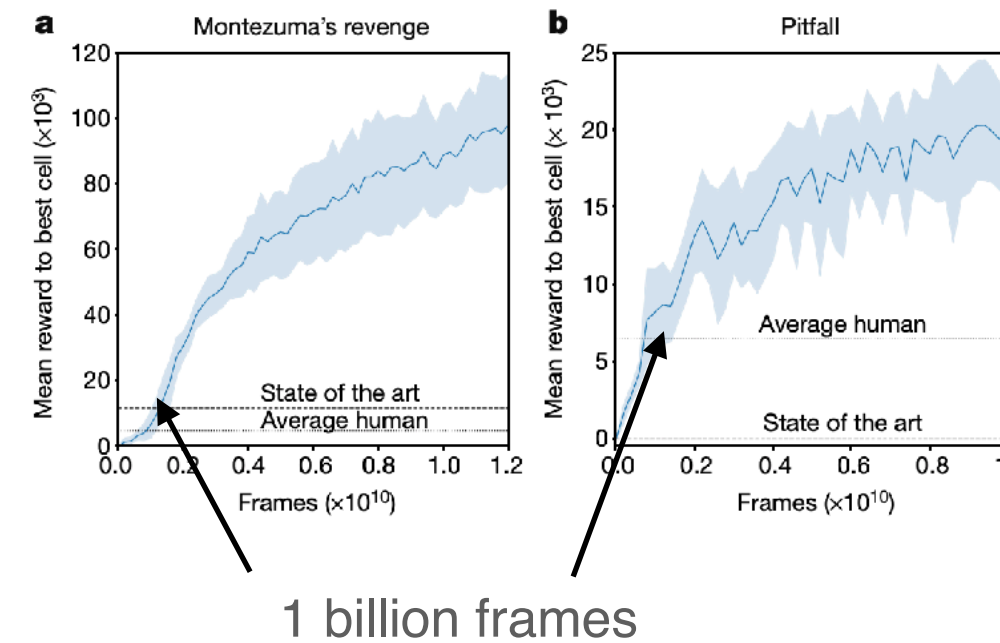


- Natural PG (NPG)[Kakade, 2001] uses a **preconditioner** to improve PG direction
- NPG is the **building block** of several **state-of-the-art** algorithms (**TRPO, PPO**)
- **Linear convergence** of NPG is established for **tabular case** [Xiao, 2022]

Context

Objective: $\arg \min_{\theta \in \mathbb{R}^d} V_{\rho}(\theta)$

- Original PG is not **sample efficient**
- Natural PG (NPG)[Kakade, 2001] uses a **preconditioner** to improve PG direction
- NPG is the **building block** of several **state-of-the-art** algorithms (**TRPO, PPO**)
- **Linear convergence** of NPG is established for **tabular case** [Xiao, 2022]



Motivations

- Extend linear convergence of NPG from tabular to **function approximation regime**.

Natural policy gradient

Natural policy gradient

- State-action cost function (a.k.a Q-function) & advantage function

Natural policy gradient

- State-action cost function (a.k.a Q-function) & advantage function

$$Q_{s,a}(\theta) := \mathbb{E}_{a_t \sim \pi_{s_t}(\theta), s_{t+1} \sim P(\cdot | s_t, a_t)} \left[\sum_{t=0}^{\infty} \gamma^t c(s_t, a_t) \mid s_0 = s, a_0 = a \right]$$

Natural policy gradient

- State-action cost function (a.k.a Q-function) & advantage function

$$Q_{s,a}(\theta) := \mathbb{E}_{a_t \sim \pi_{s_t}(\theta), s_{t+1} \sim P(\cdot | s_t, a_t)} \left[\sum_{t=0}^{\infty} \gamma^t c(s_t, a_t) \mid s_0 = s, a_0 = a \right]$$

$$A_{s,a}(\theta) := Q_{s,a}(\theta) - \mathbb{E}_{a' \sim \pi_s(\theta)} [Q_{s,a'}(\theta)]$$

Natural policy gradient

- State-action cost function (a.k.a Q-function) & advantage function

$$Q_{s,a}(\theta) := \mathbb{E}_{a_t \sim \pi_{s_t}(\theta), s_{t+1} \sim P(\cdot | s_t, a_t)} \left[\sum_{t=0}^{\infty} \gamma^t c(s_t, a_t) \mid s_0 = s, a_0 = a \right]$$

$$A_{s,a}(\theta) := Q_{s,a}(\theta) - \mathbb{E}_{a' \sim \pi_s(\theta)} [Q_{s,a'}(\theta)]$$

- Policy gradient theorem [Sutton et al., 2000]

$$\nabla_{\theta} V_{\rho}(\theta) = \frac{1}{1-\gamma} \mathbb{E}_{(s,a) \sim \mathcal{D}(\theta)} [A_{s,a}(\theta) \nabla_{\theta} \log \pi_{s,a}(\theta)]$$

Natural policy gradient

- State-action cost function (a.k.a Q-function) & advantage function

$$Q_{s,a}(\theta) := \mathbb{E}_{a_t \sim \pi_{s_t}(\theta), s_{t+1} \sim P(\cdot | s_t, a_t)} \left[\sum_{t=0}^{\infty} \gamma^t c(s_t, a_t) \mid s_0 = s, a_0 = a \right]$$

$$A_{s,a}(\theta) := Q_{s,a}(\theta) - \mathbb{E}_{a' \sim \pi_s(\theta)} [Q_{s,a'}(\theta)]$$

- Policy gradient theorem [Sutton et al., 2000]

$$\nabla_{\theta} V_{\rho}(\theta) = \frac{1}{1 - \gamma} \mathbb{E}_{(s,a) \sim \mathcal{D}(\theta)} [A_{s,a}(\theta) \nabla_{\theta} \log \pi_{s,a}(\theta)]$$

Stationary distribution of the MDP

Natural policy gradient

- State-action cost function (a.k.a Q-function) & advantage function

$$Q_{s,a}(\theta) := \mathbb{E}_{a_t \sim \pi_{s_t}(\theta), s_{t+1} \sim P(\cdot | s_t, a_t)} \left[\sum_{t=0}^{\infty} \gamma^t c(s_t, a_t) \mid s_0 = s, a_0 = a \right]$$

$$A_{s,a}(\theta) := Q_{s,a}(\theta) - \mathbb{E}_{a' \sim \pi_s(\theta)} [Q_{s,a'}(\theta)]$$

- Policy gradient theorem [Sutton et al., 2000]

$$\nabla_{\theta} V_{\rho}(\theta) = \frac{1}{1-\gamma} \mathbb{E}_{(s,a) \sim \mathcal{D}(\theta)} [A_{s,a}(\theta) \nabla_{\theta} \log \pi_{s,a}(\theta)]$$

Stationary distribution of the MDP

- Natural policy gradient

$$\theta^{(k+1)} = \theta^{(k)} - \eta_k F_{\rho}(\theta^{(k)})^{\dagger} \nabla_{\theta} V_{\rho}(\theta^{(k)})$$

Natural policy gradient

- State-action cost function (a.k.a Q-function) & advantage function

$$Q_{s,a}(\theta) := \mathbb{E}_{a_t \sim \pi_{s_t}(\theta), s_{t+1} \sim P(\cdot | s_t, a_t)} \left[\sum_{t=0}^{\infty} \gamma^t c(s_t, a_t) \mid s_0 = s, a_0 = a \right]$$

$$A_{s,a}(\theta) := Q_{s,a}(\theta) - \mathbb{E}_{a' \sim \pi_s(\theta)} [Q_{s,a'}(\theta)]$$

- Policy gradient theorem [Sutton et al., 2000]

$$\nabla_{\theta} V_{\rho}(\theta) = \frac{1}{1-\gamma} \mathbb{E}_{(s,a) \sim \mathcal{D}(\theta)} [A_{s,a}(\theta) \nabla_{\theta} \log \pi_{s,a}(\theta)]$$

Stationary distribution of the MDP

- Natural policy gradient

$$\theta^{(k+1)} = \theta^{(k)} - \eta_k F_{\rho}(\theta^{(k)})^{\dagger} \nabla_{\theta} V_{\rho}(\theta^{(k)})$$

$$F_{\rho}(\theta) = \mathbb{E}_{(s,a) \sim \mathcal{D}(\theta)} \left[\nabla_{\theta} \log \pi_{s,a}(\theta) (\nabla_{\theta} \log \pi_{s,a}(\theta))^{\top} \right] : \text{Fisher information matrix}$$

Natural policy gradient

- State-action cost function (a.k.a Q-function) & advantage function

$$Q_{s,a}(\theta) := \mathbb{E}_{a_t \sim \pi_{s_t}(\theta), s_{t+1} \sim P(\cdot | s_t, a_t)} \left[\sum_{t=0}^{\infty} \gamma^t c(s_t, a_t) \mid s_0 = s, a_0 = a \right]$$

$$A_{s,a}(\theta) := Q_{s,a}(\theta) - \mathbb{E}_{a' \sim \pi_s(\theta)} [Q_{s,a'}(\theta)]$$

- Policy gradient theorem [Sutton et al., 2000]

$$\nabla_{\theta} V_{\rho}(\theta) = \frac{1}{1-\gamma} \mathbb{E}_{(s,a) \sim \mathcal{D}(\theta)} [A_{s,a}(\theta) \nabla_{\theta} \log \pi_{s,a}(\theta)]$$

Stationary distribution of the MDP

- Natural policy gradient

$$\theta^{(k+1)} = \theta^{(k)} - \eta_k F_{\rho}(\theta^{(k)})^{\dagger} \nabla_{\theta} V_{\rho}(\theta^{(k)})$$

$$F_{\rho}(\theta) = \mathbb{E}_{(s,a) \sim \mathcal{D}(\theta)} \left[\nabla_{\theta} \log \pi_{s,a}(\theta) (\nabla_{\theta} \log \pi_{s,a}(\theta))^{\top} \right] : \text{Fisher information matrix}$$

Natural policy gradient

With log-linear policies

Log-linear policy:

$$\pi_{s,a}(\theta) = \frac{\exp \phi_{s,a}^\top \theta}{\sum_{a' \in \mathcal{A}} \exp \phi_{s,a'}^\top \theta}$$

- State-action cost function (a.k.a Q-function) & advantage function

$$Q_{s,a}(\theta) := \mathbb{E}_{a_t \sim \pi_{s_t}(\theta), s_{t+1} \sim P(\cdot | s_t, a_t)} \left[\sum_{t=0}^{\infty} \gamma^t c(s_t, a_t) \mid s_0 = s, a_0 = a \right]$$

$$A_{s,a}(\theta) := Q_{s,a}(\theta) - \mathbb{E}_{a' \sim \pi_s(\theta)} [Q_{s,a'}(\theta)]$$

- Policy gradient theorem [Sutton et al., 2000]

$$\nabla_{\theta} V_{\rho}(\theta) = \frac{1}{1-\gamma} \mathbb{E}_{(s,a) \sim \mathcal{D}(\theta)} [A_{s,a}(\theta) \nabla_{\theta} \log \pi_{s,a}(\theta)]$$

Stationary distribution of the MDP

- Natural policy gradient

$$\theta^{(k+1)} = \theta^{(k)} - \eta_k F_{\rho}(\theta^{(k)})^{\dagger} \nabla_{\theta} V_{\rho}(\theta^{(k)})$$

$$F_{\rho}(\theta) = \mathbb{E}_{(s,a) \sim \mathcal{D}(\theta)} \left[\nabla_{\theta} \log \pi_{s,a}(\theta) (\nabla_{\theta} \log \pi_{s,a}(\theta))^{\top} \right] : \text{Fisher information matrix}$$

Natural policy gradient

With log-linear policies

Log-linear policy:

$$\pi_{s,a}(\theta) = \frac{\exp \phi_{s,a}^\top \theta}{\sum_{a' \in \mathcal{A}} \exp \phi_{s,a'}^\top \theta}$$

Feature map $\phi_{s,a'} \in \mathbb{R}^d$ over $\mathcal{S} \times \mathcal{A}$

- State-action cost function (a.k.a Q-function) & advantage function

$$Q_{s,a}(\theta) := \mathbb{E}_{a_t \sim \pi_{s_t}(\theta), s_{t+1} \sim P(\cdot | s_t, a_t)} \left[\sum_{t=0}^{\infty} \gamma^t c(s_t, a_t) \mid s_0 = s, a_0 = a \right]$$

$$A_{s,a}(\theta) := Q_{s,a}(\theta) - \mathbb{E}_{a' \sim \pi_s(\theta)} [Q_{s,a'}(\theta)]$$

- Policy gradient theorem [Sutton et al., 2000]

$$\nabla_{\theta} V_{\rho}(\theta) = \frac{1}{1-\gamma} \mathbb{E}_{(s,a) \sim \mathcal{D}(\theta)} [A_{s,a}(\theta) \nabla_{\theta} \log \pi_{s,a}(\theta)]$$

Stationary distribution of the MDP

- Natural policy gradient

$$\theta^{(k+1)} = \theta^{(k)} - \eta_k F_{\rho}(\theta^{(k)})^{\dagger} \nabla_{\theta} V_{\rho}(\theta^{(k)})$$

$$F_{\rho}(\theta) = \mathbb{E}_{(s,a) \sim \mathcal{D}(\theta)} \left[\nabla_{\theta} \log \pi_{s,a}(\theta) (\nabla_{\theta} \log \pi_{s,a}(\theta))^{\top} \right] : \text{Fisher information matrix}$$

NPG with compatible function approximation

NPG with compatible function approximation

- Compatible function approximation [[Agarwal et al., 2021](#)]

$$L(w, \theta, \zeta) = \mathbb{E}_{(s,a) \sim \zeta} [(w^\top \nabla_{\theta} \log \pi_{s,a}(\theta) - A_{s,a}(\theta))^2]$$

NPG with compatible function approximation

- Compatible function approximation [[Agarwal et al., 2021](#)]

$$L(w, \theta, \zeta) = \mathbb{E}_{(s,a) \sim \zeta} [(w^\top \nabla_{\theta} \log \pi_{s,a}(\theta) - A_{s,a}(\theta))^2]$$

- NPG can be rewritten as

NPG with compatible function approximation

- Compatible function approximation [[Agarwal et al., 2021](#)]

$$L(w, \theta, \zeta) = \mathbb{E}_{(s,a) \sim \zeta} [(w^\top \nabla_{\theta} \log \pi_{s,a}(\theta) - A_{s,a}(\theta))^2]$$

- NPG can be rewritten as

$$\theta^{(k+1)} = \theta^{(k)} - \eta_k w_{\star}^{(k)}, \quad w_{\star}^{(k)} \in \arg \min_{w \in \mathbb{R}^d} L(w, \theta^{(k)}, \mathcal{D}(\theta^{(k)}))$$

NPG with compatible function approximation

- Compatible function approximation [[Agarwal et al., 2021](#)]

$$L(w, \theta, \zeta) = \mathbb{E}_{(s,a) \sim \zeta} [(w^\top \nabla_{\theta} \log \pi_{s,a}(\theta) - A_{s,a}(\theta))^2]$$

- NPG can be rewritten as

$$\theta^{(k+1)} = \theta^{(k)} - \eta_k w_{\star}^{(k)},$$

$$w_{\star}^{(k)} \in \arg \min_{w \in \mathbb{R}^d} L(w, \theta^{(k)}, \mathcal{D}(\theta^{(k)}))$$



NPG with compatible function approximation

- Compatible function approximation [Agarwal et al., 2021]

$$L(w, \theta, \zeta) = \mathbb{E}_{(s,a) \sim \zeta} \left[\underbrace{(w^\top \nabla_{\theta} \log \pi_{s,a}(\theta) - A_{s,a}(\theta))^2}_{\text{Linear approximation of the advantage function}} \right]$$

Linear approximation of the advantage function

- NPG can be rewritten as

$$\theta^{(k+1)} = \theta^{(k)} - \eta_k w_{\star}^{(k)},$$

$$w_{\star}^{(k)} \in \arg \min_{w \in \mathbb{R}^d} L(w, \theta^{(k)}, \mathcal{D}(\theta^{(k)}))$$

NPG with log-linear as policy mirror descent

NPG with log-linear as policy mirror descent

Log-linear policy:

$$\pi_{s,a}(\theta) = \frac{\exp \phi_{s,a}^\top \theta}{\sum_{a' \in \mathcal{A}} \exp \phi_{s,a'}^\top \theta}$$

NPG with log-linear as policy mirror descent

- NPG with log-linear can also be written as

Log-linear policy:

$$\pi_{s,a}(\theta) = \frac{\exp \phi_{s,a}^\top \theta}{\sum_{a' \in \mathcal{A}} \exp \phi_{s,a'}^\top \theta}$$

NPG with log-linear as policy mirror descent

Log-linear policy:

$$\pi_{s,a}(\theta) = \frac{\exp \phi_{s,a}^\top \theta}{\sum_{a' \in \mathcal{A}} \exp \phi_{s,a'}^\top \theta}$$

- NPG with log-linear can also be written as

$$\pi_s(\theta^{(k+1)}) = \arg \min_{p \in \Delta(\mathcal{A})} \left\{ \eta_k \langle \bar{\Phi}_s^{(k)} w_\star^{(k)}, p \rangle + \text{KL}(p, \pi_s(\theta^{(k)})) \right\}$$

NPG with log-linear as policy mirror descent

Log-linear policy:

$$\pi_{s,a}(\theta) = \frac{\exp \phi_{s,a}^\top \theta}{\sum_{a' \in \mathcal{A}} \exp \phi_{s,a'}^\top \theta}$$

- NPG with log-linear can also be written as

$$\pi_s(\theta^{(k+1)}) = \arg \min_{p \in \Delta(\mathcal{A})} \left\{ \eta_k \langle \bar{\Phi}_s^{(k)} w_\star^{(k)}, p \rangle + \text{KL}(p, \pi_s(\theta^{(k)})) \right\} \rightarrow \text{Policy mirror descent}$$

NPG with log-linear as policy mirror descent

Log-linear policy:

$$\pi_{s,a}(\theta) = \frac{\exp \phi_{s,a}^\top \theta}{\sum_{a' \in \mathcal{A}} \exp \phi_{s,a'}^\top \theta}$$

- NPG with log-linear can also be written as

$$\pi_s(\theta^{(k+1)}) = \arg \min_{p \in \Delta(\mathcal{A})} \left\{ \eta_k \langle \bar{\Phi}_s^{(k)} w_\star^{(k)}, p \rangle + \text{KL}(p, \pi_s(\theta^{(k)})) \right\} \rightarrow \text{Policy mirror descent}$$

$\bar{\Phi}_s^{(k)} \in \mathbb{R}^{|\mathcal{A}| \times d}$ is a matrix whose rows consist of the *centered feature maps* $\bar{\phi}_{s,a}(\theta^{(k)})$

NPG with log-linear as policy mirror descent

Log-linear policy:

$$\pi_{s,a}(\theta) = \frac{\exp \phi_{s,a}^\top \theta}{\sum_{a' \in \mathcal{A}} \exp \phi_{s,a'}^\top \theta}$$

- NPG with log-linear can also be written as

$$\pi_s(\theta^{(k+1)}) = \arg \min_{p \in \Delta(\mathcal{A})} \left\{ \eta_k \langle \bar{\Phi}_s^{(k)} w_\star^{(k)}, p \rangle + \text{KL}(p, \pi_s(\theta^{(k)})) \right\} \rightarrow \text{Policy mirror descent}$$

$\bar{\Phi}_s^{(k)} \in \mathbb{R}^{|\mathcal{A}| \times d}$ is a matrix whose rows consist of the *centered feature maps* $\bar{\phi}_{s,a}(\theta^{(k)})$

$$\bar{\phi}_{s,a}(\theta^{(k)}) := \nabla_{\theta} \log \pi_{s,a}(\theta^{(k)}) = \phi_{s,a} - \mathbb{E}_{a' \sim \pi_s(\theta^{(k)})} [\phi_{s,a'}]$$

NPG with log-linear as policy mirror descent

Log-linear policy:

$$\pi_{s,a}(\theta) = \frac{\exp \phi_{s,a}^\top \theta}{\sum_{a' \in \mathcal{A}} \exp \phi_{s,a'}^\top \theta}$$

- NPG with log-linear can also be written as

$$\pi_s(\theta^{(k+1)}) = \arg \min_{p \in \Delta(\mathcal{A})} \left\{ \eta_k \langle \bar{\Phi}_s^{(k)} w_\star^{(k)}, p \rangle + \text{KL}(p, \pi_s(\theta^{(k)})) \right\} \rightarrow \text{Policy mirror descent}$$

$\bar{\Phi}_s^{(k)} \in \mathbb{R}^{|\mathcal{A}| \times d}$ is a matrix whose rows consist of the *centered feature maps* $\bar{\phi}_{s,a}(\theta^{(k)})$

$$\bar{\phi}_{s,a}(\theta^{(k)}) := \nabla_{\theta} \log \pi_{s,a}(\theta^{(k)}) = \phi_{s,a} - \mathbb{E}_{a' \sim \pi_s(\theta^{(k)})} [\phi_{s,a'}]$$

$\text{KL}(p, q) = \sum_{a \in \mathcal{A}} p_a \log(p_a/q_a)$ is the Kullback-Leibler (KL) divergence for $p, q \in \Delta(\mathcal{A})$

NPG with log-linear as policy mirror descent

Log-linear policy:

$$\pi_{s,a}(\theta) = \frac{\exp \phi_{s,a}^\top \theta}{\sum_{a' \in \mathcal{A}} \exp \phi_{s,a'}^\top \theta}$$

- NPG with log-linear can also be written as

$$\pi_s(\theta^{(k+1)}) = \arg \min_{p \in \Delta(\mathcal{A})} \left\{ \eta_k \langle \bar{\Phi}_s^{(k)} w_\star^{(k)}, p \rangle + \text{KL}(p, \pi_s(\theta^{(k)})) \right\} \rightarrow \text{Policy mirror descent}$$

$\bar{\Phi}_s^{(k)} \in \mathbb{R}^{|\mathcal{A}| \times d}$ is a matrix whose rows consist of the *centered feature maps* $\bar{\phi}_{s,a}(\theta^{(k)})$

$$\bar{\phi}_{s,a}(\theta^{(k)}) := \nabla_{\theta} \log \pi_{s,a}(\theta^{(k)}) = \phi_{s,a} - \mathbb{E}_{a' \sim \pi_s(\theta^{(k)})} [\phi_{s,a'}]$$

$\text{KL}(p, q) = \sum_{a \in \mathcal{A}} p_a \log(p_a/q_a)$ is the Kullback-Leibler (KL) divergence for $p, q \in \Delta(\mathcal{A})$

- Connection with Policy Iteration

NPG with log-linear as policy mirror descent

Log-linear policy:

$$\pi_{s,a}(\theta) = \frac{\exp \phi_{s,a}^\top \theta}{\sum_{a' \in \mathcal{A}} \exp \phi_{s,a'}^\top \theta}$$

- NPG with log-linear can also be written as

$$\pi_s(\theta^{(k+1)}) = \arg \min_{p \in \Delta(\mathcal{A})} \left\{ \eta_k \langle \bar{\Phi}_s^{(k)} w_\star^{(k)}, p \rangle + \text{KL}(p, \pi_s(\theta^{(k)})) \right\} \rightarrow \text{Policy mirror descent}$$

$\bar{\Phi}_s^{(k)} \in \mathbb{R}^{|\mathcal{A}| \times d}$ is a matrix whose rows consist of the *centered feature maps* $\bar{\phi}_{s,a}(\theta^{(k)})$

$$\bar{\phi}_{s,a}(\theta^{(k)}) := \nabla_{\theta} \log \pi_{s,a}(\theta^{(k)}) = \phi_{s,a} - \mathbb{E}_{a' \sim \pi_s(\theta^{(k)})} [\phi_{s,a'}]$$

$\text{KL}(p, q) = \sum_{a \in \mathcal{A}} p_a \log(p_a/q_a)$ is the Kullback-Leibler (KL) divergence for $p, q \in \Delta(\mathcal{A})$

- Connection with Policy Iteration

$$\pi_s(\theta^{(k+1)}) = \arg \min_{p \in \Delta(\mathcal{A})} \left\{ \eta_k \langle A_s(\theta^{(k)}), p \rangle \right\} \quad \text{with } A_s(\theta^{(k)}) := [A_{s,a}(\theta^{(k)})]_a \in \mathbb{R}^{|\mathcal{A}|}$$

NPG with log-linear as policy mirror descent

Log-linear policy:

$$\pi_{s,a}(\theta) = \frac{\exp \phi_{s,a}^\top \theta}{\sum_{a' \in \mathcal{A}} \exp \phi_{s,a'}^\top \theta}$$

- NPG with log-linear can also be written as

$$\pi_s(\theta^{(k+1)}) = \arg \min_{p \in \Delta(\mathcal{A})} \left\{ \eta_k \langle \bar{\Phi}_s^{(k)} w_\star^{(k)}, p \rangle + \text{KL}(p, \pi_s(\theta^{(k)})) \right\} \rightarrow \text{Policy mirror descent}$$

$\bar{\Phi}_s^{(k)} \in \mathbb{R}^{|\mathcal{A}| \times d}$ is a matrix whose rows consist of the *centered feature maps* $\bar{\phi}_{s,a}(\theta^{(k)})$

$$\bar{\phi}_{s,a}(\theta^{(k)}) := \nabla_{\theta} \log \pi_{s,a}(\theta^{(k)}) = \phi_{s,a} - \mathbb{E}_{a' \sim \pi_s(\theta^{(k)})} [\phi_{s,a'}] \quad \text{Regularization}$$

$\text{KL}(p, q) = \sum_{a \in \mathcal{A}} p_a \log(p_a/q_a)$ is the Kullback-Leibler (KL) divergence for $p, q \in \Delta(\mathcal{A})$

- Connection with Policy Iteration

$$\pi_s(\theta^{(k+1)}) = \arg \min_{p \in \Delta(\mathcal{A})} \left\{ \eta_k \langle A_s(\theta^{(k)}), p \rangle \right\} \quad \text{with } A_s(\theta^{(k)}) := [A_{s,a}(\theta^{(k)})]_a \in \mathbb{R}^{|\mathcal{A}|}$$

NPG with log-linear as policy mirror descent

Log-linear policy:

$$\pi_{s,a}(\theta) = \frac{\exp \phi_{s,a}^\top \theta}{\sum_{a' \in \mathcal{A}} \exp \phi_{s,a'}^\top \theta}$$

- NPG with log-linear can also be written as

$$\pi_s(\theta^{(k+1)}) = \arg \min_{p \in \Delta(\mathcal{A})} \left\{ \eta_k \langle \bar{\Phi}_s^{(k)} w_\star^{(k)}, p \rangle + \text{KL}(p, \pi_s(\theta^{(k)})) \right\} \rightarrow \text{Policy mirror descent}$$

$\bar{\Phi}_s^{(k)} \in \mathbb{R}^{|\mathcal{A}| \times d}$ is a matrix whose rows consist of the *centered feature maps* $\bar{\phi}_{s,a}(\theta^{(k)})$

$$\bar{\phi}_{s,a}(\theta^{(k)}) := \nabla_{\theta} \log \pi_{s,a}(\theta^{(k)}) = \phi_{s,a} - \mathbb{E}_{a' \sim \pi_s(\theta^{(k)})} [\phi_{s,a'}] \quad \text{Regularization}$$

$\text{KL}(p, q) = \sum_{a \in \mathcal{A}} p_a \log(p_a/q_a)$ is the Kullback-Leibler (KL) divergence for $p, q \in \Delta(\mathcal{A})$

- Connection with Policy Iteration

$$\pi_s(\theta^{(k+1)}) = \arg \min_{p \in \Delta(\mathcal{A})} \left\{ \eta_k \langle A_s(\theta^{(k)}), p \rangle \right\} \quad \text{with } A_s(\theta^{(k)}) := [A_{s,a}(\theta^{(k)})]_a \in \mathbb{R}^{|\mathcal{A}|}$$

Linear approximation

Two key ingredients

Two key ingredients

- Performance difference lemma [[Kakade and Langford, 2002](#)]: For any policy $\pi(\theta), \pi(\theta')$,

Two key ingredients

- Performance difference lemma [[Kakade and Langford, 2002](#)]: For any policy $\pi(\theta), \pi(\theta')$,

$$V_{\rho}(\theta) - V_{\rho}(\theta') = \frac{1}{1 - \gamma} \mathbb{E}_{(s,a) \sim \mathcal{D}(\theta)} [A_{s,a}(\theta')]$$

Two key ingredients

- Performance difference lemma [[Kakade and Langford, 2002](#)]: For any policy $\pi(\theta), \pi(\theta')$,

$$V_{\rho}(\theta) - V_{\rho}(\theta') = \frac{1}{1 - \gamma} \mathbb{E}_{(s,a) \sim \mathcal{D}(\theta)} [A_{s,a}(\theta')]$$

- Three-point descent lemma [[Chen and Teboulle, 1993](#)]:

Two key ingredients

- Performance difference lemma [[Kakade and Langford, 2002](#)]: For any policy $\pi(\theta), \pi(\theta')$,

$$V_\rho(\theta) - V_\rho(\theta') = \frac{1}{1 - \gamma} \mathbb{E}_{(s,a) \sim \mathcal{D}(\theta)} [A_{s,a}(\theta')]$$

- Three-point descent lemma [[Chen and Teboulle, 1993](#)]:

Given a convex set \mathcal{C} and a convex function $f : \mathcal{C} \rightarrow \mathbb{R}$, for any $q \in \mathcal{C}$, let

Two key ingredients

- Performance difference lemma [[Kakade and Langford, 2002](#)]: For any policy $\pi(\theta), \pi(\theta')$,

$$V_\rho(\theta) - V_\rho(\theta') = \frac{1}{1-\gamma} \mathbb{E}_{(s,a) \sim \mathcal{D}(\theta)} [A_{s,a}(\theta')]$$

- Three-point descent lemma [[Chen and Teboulle, 1993](#)]:

Given a convex set \mathcal{C} and a convex function $f : \mathcal{C} \rightarrow \mathbb{R}$, for any $q \in \mathcal{C}$, let

$$p^+ \in \arg \min_{p \in \mathcal{C}} f(p) + \text{KL}(p, q)$$

Two key ingredients

- Performance difference lemma [[Kakade and Langford, 2002](#)]: For any policy $\pi(\theta), \pi(\theta')$,

$$V_\rho(\theta) - V_\rho(\theta') = \frac{1}{1 - \gamma} \mathbb{E}_{(s,a) \sim \mathcal{D}(\theta)} [A_{s,a}(\theta')]$$

- Three-point descent lemma [[Chen and Teboulle, 1993](#)]:

Given a convex set \mathcal{C} and a convex function $f : \mathcal{C} \rightarrow \mathbb{R}$, for any $q \in \mathcal{C}$, let

$$p^+ \in \arg \min_{p \in \mathcal{C}} f(p) + \text{KL}(p, q)$$

$$\implies f(p^+) + \text{KL}(p^+, q) \leq f(p) + \text{KL}(p, q) - \text{KL}(p, p^+), \quad \forall p \in \mathcal{C}$$

Two key ingredients

- Performance difference lemma [[Kakade and Langford, 2002](#)]: For any policy $\pi(\theta), \pi(\theta')$,

$$V_\rho(\theta) - V_\rho(\theta') = \frac{1}{1-\gamma} \mathbb{E}_{(s,a) \sim \mathcal{D}(\theta)} [A_{s,a}(\theta')]$$

- Three-point descent lemma [[Chen and Teboulle, 1993](#)]:

Given a convex set \mathcal{C} and a convex function $f: \mathcal{C} \rightarrow \mathbb{R}$, for any $q \in \mathcal{C}$, let

$$p^+ \in \arg \min_{p \in \mathcal{C}} f(p) + \text{KL}(p, q)$$

$$\implies f(p^+) + \text{KL}(p^+, q) \leq f(p) + \text{KL}(p, q) - \text{KL}(p, p^+), \quad \forall p \in \mathcal{C}$$

$$\text{💡 } \pi_s(\theta^{(k+1)}) = \arg \min_{p \in \Delta(\mathcal{A})} \left\{ \eta_k \langle \bar{\Phi}_s^{(k)} w_\star^{(k)}, p \rangle + \text{KL}(p, \pi_s(\theta^{(k)})) \right\}$$

Two key ingredients

- Performance difference lemma [[Kakade and Langford, 2002](#)]: For any policy $\pi(\theta), \pi(\theta')$,

$$V_\rho(\theta) - V_\rho(\theta') = \frac{1}{1-\gamma} \mathbb{E}_{(s,a) \sim \mathcal{D}(\theta)} [A_{s,a}(\theta')]$$

- Three-point descent lemma [[Chen and Teboulle, 1993](#)]:

Given a convex set \mathcal{C} and a convex function $f: \mathcal{C} \rightarrow \mathbb{R}$, for any $q \in \mathcal{C}$, let

$$p^+ \in \arg \min_{p \in \mathcal{C}} f(p) + \text{KL}(p, q)$$

$$\implies f(p^+) + \text{KL}(p^+, q) \leq f(p) + \text{KL}(p, q) - \text{KL}(p, p^+), \quad \forall p \in \mathcal{C}$$

$$\text{💡 } \pi_s(\theta^{(k+1)}) = \arg \min_{p \in \Delta(\mathcal{A})} \left\{ \eta_k \langle \bar{\Phi}_s^{(k)} w_\star^{(k)}, p \rangle + \text{KL}(p, \pi_s(\theta^{(k)})) \right\}$$

Two key ingredients

- Performance difference lemma [[Kakade and Langford, 2002](#)]: For any policy $\pi(\theta), \pi(\theta')$,

$$V_\rho(\theta) - V_\rho(\theta') = \frac{1}{1-\gamma} \mathbb{E}_{(s,a) \sim \mathcal{D}(\theta)} [A_{s,a}(\theta')]$$

- Three-point descent lemma [[Chen and Teboulle, 1993](#)]:

Given a convex set \mathcal{C} and a convex function $f: \mathcal{C} \rightarrow \mathbb{R}$, for any $q \in \mathcal{C}$, let

$$p^+ \in \arg \min_{p \in \mathcal{C}} f(p) + \text{KL}(p, q)$$

$$\implies f(p^+) + \text{KL}(p^+, q) \leq f(p) + \text{KL}(p, q) - \text{KL}(p, p^+), \quad \forall p \in \mathcal{C}$$

$$\text{💡 } \pi_s(\theta^{(k+1)}) = \arg \min_{p \in \Delta(\mathcal{A})} \left\{ \eta_k \langle \bar{\Phi}_s^{(k)} w_\star^{(k)}, p \rangle + \text{KL}(p, \pi_s(\theta^{(k)})) \right\}$$

Two key ingredients

- Performance difference lemma [[Kakade and Langford, 2002](#)]: For any policy $\pi(\theta), \pi(\theta')$,

$$V_\rho(\theta) - V_\rho(\theta') = \frac{1}{1-\gamma} \mathbb{E}_{(s,a) \sim \mathcal{D}(\theta)} [A_{s,a}(\theta')]$$

- Three-point descent lemma [[Chen and Teboulle, 1993](#)]:

Given a convex set \mathcal{C} and a convex function $f: \mathcal{C} \rightarrow \mathbb{R}$, for any $q \in \mathcal{C}$, let

$$p^+ \in \arg \min_{p \in \mathcal{C}} f(p) + \text{KL}(p, q)$$

$$\implies f(p^+) + \text{KL}(p^+, q) \leq f(p) + \text{KL}(p, q) - \text{KL}(p, p^+), \quad \forall p \in \mathcal{C}$$

$$\underbrace{\text{lightbulb}}_{\pi_s(\theta^{(k+1)})} = \arg \min_{p \in \Delta(\mathcal{A})} \left\{ \eta_k \langle \bar{\Phi}_s^{(k)} w_\star^{(k)}, p \rangle + \text{KL}(p, \pi_s(\theta^{(k)})) \right\}$$

Two key ingredients

- Performance difference lemma [[Kakade and Langford, 2002](#)]: For any policy $\pi(\theta), \pi(\theta')$,

$$V_\rho(\theta) - V_\rho(\theta') = \frac{1}{1-\gamma} \mathbb{E}_{(s,a) \sim \mathcal{D}(\theta)} [A_{s,a}(\theta')]$$

- Three-point descent lemma [[Chen and Teboulle, 1993](#)]:

Given a convex set \mathcal{C} and a convex function $f: \mathcal{C} \rightarrow \mathbb{R}$, for any $q \in \mathcal{C}$, let

$$p^+ \in \arg \min_{p \in \mathcal{C}} f(p) + \text{KL}(p, q)$$

$$\implies f(p^+) + \text{KL}(p^+, q) \leq f(p) + \text{KL}(p, q) - \text{KL}(p, p^+), \quad \forall p \in \mathcal{C}$$

$$\underbrace{\text{lightbulb}}_{\pi_s(\theta^{(k+1)})} = \arg \min_{p \in \Delta(\mathcal{A})} \left\{ \underbrace{\eta_k \langle \bar{\Phi}_s^{(k)} w_\star^{(k)}, p \rangle}_{\text{inner product}} + \underbrace{\text{KL}(p, \pi_s(\theta^{(k)}))}_{\text{KL divergence}} \right\}$$

Convergence theory

(see paper for technique details and additional properties)

Convergence theory

(see paper for technique details and additional properties)

- Three-point descent lemma :

For any $p \in \Delta(\mathcal{A})$,

Convergence theory

(see paper for technique details and additional properties)

- Three-point descent lemma :

For any $p \in \Delta(\mathcal{A})$,

$$\begin{aligned} & \eta_k \langle \bar{\Phi}_s^{(k)} w_{\star}^{(k)}, \pi_s(\theta^{(k+1)}) \rangle + \text{KL}(\pi_s(\theta^{(k+1)}), \pi_s(\theta^{(k)})) \\ & \leq \eta_k \langle \bar{\Phi}_s^{(k)} w_{\star}^{(k)}, p \rangle + \text{KL}(p, \pi_s(\theta^{(k)})) - \text{KL}(p, \pi_s(\theta^{(k+1)})) \end{aligned}$$

Convergence theory

(see paper for technique details and additional properties)

- Three-point descent lemma :

For any $p \in \Delta(\mathcal{A})$,

$$\begin{aligned} & \eta_k \langle \bar{\Phi}_s^{(k)} w_{\star}^{(k)}, \pi_s(\theta^{(k+1)}) \rangle + \text{KL}(\pi_s(\theta^{(k+1)}), \pi_s(\theta^{(k)})) \\ & \leq \eta_k \langle \bar{\Phi}_s^{(k)} w_{\star}^{(k)}, p \rangle + \text{KL}(p, \pi_s(\theta^{(k)})) - \text{KL}(p, \pi_s(\theta^{(k+1)})) \end{aligned}$$

One can let $p = \pi_s(\theta^{(k)})$ or be the optimal policy to derive a **telescoping sum**

Convergence theory

(see paper for technique details and additional properties)

- Three-point descent lemma :

For any $p \in \Delta(\mathcal{A})$,

$$\begin{aligned} & \eta_k \langle \bar{\Phi}_s^{(k)} w_{\star}^{(k)}, \pi_s(\theta^{(k+1)}) \rangle + \text{KL}(\pi_s(\theta^{(k+1)}), \pi_s(\theta^{(k)})) \\ & \leq \eta_k \langle \bar{\Phi}_s^{(k)} w_{\star}^{(k)}, p \rangle + \text{KL}(p, \pi_s(\theta^{(k)})) - \text{KL}(p, \pi_s(\theta^{(k+1)})) \end{aligned}$$

One can let $p = \pi_s(\theta^{(k)})$ or be the optimal policy to derive a **telescoping sum**

- **Linear convergence** to the **global optimum** by increasing step size by $1/\gamma$

Convergence theory

(see paper for technique details and additional properties)

- Three-point descent lemma :

For any $p \in \Delta(\mathcal{A})$,

$$\begin{aligned} & \eta_k \langle \bar{\Phi}_s^{(k)} w_{\star}^{(k)}, \pi_s(\theta^{(k+1)}) \rangle + \text{KL}(\pi_s(\theta^{(k+1)}), \pi_s(\theta^{(k)})) \\ & \leq \eta_k \langle \bar{\Phi}_s^{(k)} w_{\star}^{(k)}, p \rangle + \text{KL}(p, \pi_s(\theta^{(k)})) - \text{KL}(p, \pi_s(\theta^{(k+1)})) \end{aligned}$$

One can let $p = \pi_s(\theta^{(k)})$ or be the optimal policy to derive a **telescoping sum**

- **Linear convergence** to the **global optimum** by increasing step size by $1/\gamma$

$$\pi_s(\theta^{(k+1)}) = \arg \min_{p \in \Delta(\mathcal{A})} \left\{ \eta_k \langle \bar{\Phi}_s^{(k)} w_{\star}^{(k)}, p \rangle + \text{KL}(p, \pi_s(\theta^{(k)})) \right\}$$

Convergence theory

(see paper for technique details and additional properties)

- Three-point descent lemma :

For any $p \in \Delta(\mathcal{A})$,

$$\begin{aligned} & \eta_k \langle \bar{\Phi}_s^{(k)} w_{\star}^{(k)}, \pi_s(\theta^{(k+1)}) \rangle + \text{KL}(\pi_s(\theta^{(k+1)}), \pi_s(\theta^{(k)})) \\ & \leq \eta_k \langle \bar{\Phi}_s^{(k)} w_{\star}^{(k)}, p \rangle + \text{KL}(p, \pi_s(\theta^{(k)})) - \text{KL}(p, \pi_s(\theta^{(k+1)})) \end{aligned}$$

One can let $p = \pi_s(\theta^{(k)})$ or be the optimal policy to derive a **telescoping sum**

- **Linear convergence** to the **global optimum** by increasing step size by $1/\gamma$

$$\pi_s(\theta^{(k+1)}) = \arg \min_{p \in \Delta(\mathcal{A})} \left\{ \eta_k \langle \bar{\Phi}_s^{(k)} w_{\star}^{(k)}, p \rangle + \text{KL}(p, \pi_s(\theta^{(k)})) \right\} \quad \eta_k \longrightarrow \infty$$

Convergence theory

(see paper for technique details and additional properties)

- Three-point descent lemma :

For any $p \in \Delta(\mathcal{A})$,

$$\begin{aligned} & \eta_k \langle \bar{\Phi}_s^{(k)} w_{\star}^{(k)}, \pi_s(\theta^{(k+1)}) \rangle + \text{KL}(\pi_s(\theta^{(k+1)}), \pi_s(\theta^{(k)})) \\ & \leq \eta_k \langle \bar{\Phi}_s^{(k)} w_{\star}^{(k)}, p \rangle + \text{KL}(p, \pi_s(\theta^{(k)})) - \text{KL}(p, \pi_s(\theta^{(k+1)})) \end{aligned}$$

One can let $p = \pi_s(\theta^{(k)})$ or be the optimal policy to derive a **telescoping sum**

- **Linear convergence** to the **global optimum** by increasing step size by $1/\gamma$

$$\pi_s(\theta^{(k+1)}) = \arg \min_{p \in \Delta(\mathcal{A})} \left\{ \eta_k \langle \bar{\Phi}_s^{(k)} w_{\star}^{(k)}, p \rangle + \text{KL}(p, \pi_s(\theta^{(k)})) \right\} \quad \eta_k \longrightarrow \infty$$

Convergence theory

(see paper for technique details and additional properties)

- Three-point descent lemma :

For any $p \in \Delta(\mathcal{A})$,

$$\begin{aligned} & \eta_k \langle \bar{\Phi}_s^{(k)} w_{\star}^{(k)}, \pi_s(\theta^{(k+1)}) \rangle + \text{KL}(\pi_s(\theta^{(k+1)}), \pi_s(\theta^{(k)})) \\ & \leq \eta_k \langle \bar{\Phi}_s^{(k)} w_{\star}^{(k)}, p \rangle + \text{KL}(p, \pi_s(\theta^{(k)})) - \text{KL}(p, \pi_s(\theta^{(k+1)})) \end{aligned}$$

One can let $p = \pi_s(\theta^{(k)})$ or be the optimal policy to derive a **telescoping sum**

- **Linear convergence** to the **global optimum** by increasing step size by $1/\gamma$

$$\pi_s(\theta^{(k+1)}) = \arg \min_{p \in \Delta(\mathcal{A})} \left\{ \eta_k \langle \bar{\Phi}_s^{(k)} w_{\star}^{(k)}, p \rangle + \text{KL}(p, \pi_s(\theta^{(k)})) \right\} \quad \eta_k \longrightarrow \infty$$

Convergence theory

(see paper for technique details and additional properties)

- Three-point descent lemma :

For any $p \in \Delta(\mathcal{A})$,

$$\begin{aligned} & \eta_k \langle \bar{\Phi}_s^{(k)} w_{\star}^{(k)}, \pi_s(\theta^{(k+1)}) \rangle + \text{KL}(\pi_s(\theta^{(k+1)}), \pi_s(\theta^{(k)})) \\ & \leq \eta_k \langle \bar{\Phi}_s^{(k)} w_{\star}^{(k)}, p \rangle + \text{KL}(p, \pi_s(\theta^{(k)})) - \text{KL}(p, \pi_s(\theta^{(k+1)})) \end{aligned}$$

One can let $p = \pi_s(\theta^{(k)})$ or be the optimal policy to derive a **telescoping sum**

- **Linear convergence** to the **global optimum** by increasing step size by $1/\gamma$

$$\pi_s(\theta^{(k+1)}) = \arg \min_{p \in \Delta(\mathcal{A})} \left\{ \eta_k \langle \bar{\Phi}_s^{(k)} w_{\star}^{(k)}, p \rangle + \text{KL}(p, \pi_s(\theta^{(k)})) \right\} \quad \eta_k \longrightarrow \infty$$

 Behave more and more like policy iteration

Convergence theory 2

Convergence theory 2

- Consequently, we obtain an $\tilde{O}(\epsilon^{-2})$ sample complexity for NPG

Convergence theory 2

- Consequently, we obtain an $\tilde{O}(\epsilon^{-2})$ sample complexity for NPG
- Similar linear convergence and $\tilde{O}(\epsilon^{-2})$ sample complexity results are also established for **Q-NPG**

Convergence theory 2

- Consequently, we obtain an $\tilde{O}(\epsilon^{-2})$ sample complexity for NPG
- Similar linear convergence and $\tilde{O}(\epsilon^{-2})$ sample complexity results are also established for **Q-NPG**
- Sublinear convergence $1/K$ for both NPG and Q-NPG with arbitrary large constant step size

Discussion & Conclusion

Discussion & Conclusion

- We show that NPG (and Q-NPG) with log-linear policies enjoy **linear convergence rates** and $O(1/\varepsilon^2)$ sample complexities using a simple, **non-adaptive geometrically increasing** step size.

Discussion & Conclusion

- We show that NPG (and Q-NPG) with log-linear policies enjoy **linear convergence rates** and $O(1/\varepsilon^2)$ sample complexities using a simple, **non-adaptive geometrically increasing** step size.
- The linear convergence analysis of NPG with log-linear policy can be **extended to general parametrization** [A Novel Framework for Policy Mirror Descent with General Parametrization and Linear Convergence, Carlo Alfano, Rui Yuan, and Patrick Rebeschini, 2023].

Thank you !



References

- ▶ Vijay Konda and John Tsitsiklis. Actor-critic algorithms. In *Advances in Neural Information Processing Systems*, volume 12. MIT Press, 2000.
- ▶ Sham M Kakade. A natural policy gradient. In *Advances in Neural Information Processing Systems*, volume 14. MIT Press, 2001.
- ▶ John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region policy optimization. In *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pages 1889–1897, Lille, France, 07–09 Jul 2015.
- ▶ John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms, 2017.
- ▶ Matteo Papini, Damiano Binaghi, Giuseppe Canonaco, Matteo Pirota, and Marcello Restelli. Stochastic variance-reduced policy gradient. In *Proceedings of the 35th International Conference on Machine Learning*, volume 80, pages 4026–4035. PMLR, 2018.
- ▶ Zebang Shen, Alejandro Ribeiro, Hamed Hassani, Hui Qian, and Chao Mi. Hessian aided policy gradient. In *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 5729– 5738. PMLR, 09–15 Jun 2019
- ▶ Pan Xu, Felicia Gao, and Quanquan Gu. Sample efficient policy gradient methods with recursive variance reduction. In *International Conference on Learning Representations*, 2020.
- ▶ Feihu Huang, Shangqian Gao, Jian Pei, and Heng Huang. Momentum-based policy gradient methods, 2020.
- ▶ Lin Xiao. On the convergence rates of policy gradient methods. *Journal of Machine Learning Research*, 23(282):1–36, 2022.
- ▶ Alekh Agarwal, Sham M. Kakade, Jason D. Lee, and Gaurav Mahajan. On the theory of policy gradient methods: Optimality, approximation, and distribution shift. *Journal of Machine Learning Research* 22.98, pp. 1–76, 2021.
- ▶ Richard S Sutton, David A. McAllester, Satinder P. Singh, and Yishay Mansour. Policy gradient methods for reinforcement learning with function approximation. In *Advances in Neural Information Processing Systems* 12, pages 1057–1063. MIT Press, 2000.
- ▶ Gong Chen and Marc Teboulle. Convergence analysis of a proximal-like minimization algorithm using bregman functions. *SIAM Journal on Optimization*, 3(3):538–543, 1993.
- ▶ Sham Kakade and John Langford. Approximately optimal approximate reinforcement learning. In *Proceedings of 19th International Conference on Machine Learning*, pages 267–274, 2002.
- ▶ Carlo Alfano, Rui Yuan, and Patrick Rebeschini. A Novel Framework for Policy Mirror Descent with General Parametrization and Linear Convergence, 2023.

PG method: $\theta^{(k+1)} = \theta^{(k)} - \eta_k \nabla_{\theta} V_{\rho}(\theta^{(k)})$

Policy gradient theorem (PGT) [Sutton et al., 2000]

PG method: $\theta^{(k+1)} = \theta^{(k)} - \eta_k \nabla_{\theta} V_{\rho}(\theta^{(k)})$

Policy gradient theorem (PGT) [Sutton et al., 2000]

PG method: $\theta^{(k+1)} = \theta^{(k)} - \eta_k \nabla_{\theta} V_{\rho}(\theta^{(k)})$

- State-action cost function (a.k.a Q-function) & advantage function

Policy gradient theorem (PGT) [Sutton et al., 2000]

PG method: $\theta^{(k+1)} = \theta^{(k)} - \eta_k \nabla_{\theta} V_{\rho}(\theta^{(k)})$

- State-action cost function (a.k.a Q-function) & advantage function

$$Q_{s,a}(\theta) := \mathbb{E}_{a_t \sim \pi_{s_t}(\theta), s_{t+1} \sim P(\cdot | s_t, a_t)} \left[\sum_{t=0}^{\infty} \gamma^t c(s_t, a_t) \mid s_0 = s, a_0 = a \right]$$

Policy gradient theorem (PGT) [Sutton et al., 2000]

PG method: $\theta^{(k+1)} = \theta^{(k)} - \eta_k \nabla_{\theta} V_{\rho}(\theta^{(k)})$

- State-action cost function (a.k.a Q-function) & advantage function

$$Q_{s,a}(\theta) := \mathbb{E}_{a_t \sim \pi_{s_t}(\theta), s_{t+1} \sim P(\cdot | s_t, a_t)} \left[\sum_{t=0}^{\infty} \gamma^t c(s_t, a_t) \mid s_0 = s, a_0 = a \right]$$

$$\begin{aligned} V_{\rho}(\theta) &= \mathbb{E}_{s \sim \rho} [V_s(\theta)] \\ &= \mathbb{E}_{s \sim \rho, a \sim \pi_s(\theta)} [Q_{s,a}(\theta)] \end{aligned}$$

Policy gradient theorem (PGT) [Sutton et al., 2000]

PG method: $\theta^{(k+1)} = \theta^{(k)} - \eta_k \nabla_{\theta} V_{\rho}(\theta^{(k)})$

- State-action cost function (a.k.a Q-function) & advantage function

$$Q_{s,a}(\theta) := \mathbb{E}_{a_t \sim \pi_{s_t}(\theta), s_{t+1} \sim P(\cdot | s_t, a_t)} \left[\sum_{t=0}^{\infty} \gamma^t c(s_t, a_t) \mid s_0 = s, a_0 = a \right]$$

$$A_{s,a}(\theta) := Q_{s,a}(\theta) - \mathbb{E}_{a' \sim \pi_s(\theta)} [Q_{s,a'}(\theta)]$$

$$\begin{aligned} V_{\rho}(\theta) &= \mathbb{E}_{s \sim \rho} [V_s(\theta)] \\ &= \mathbb{E}_{s \sim \rho, a \sim \pi_s(\theta)} [Q_{s,a}(\theta)] \end{aligned}$$

Policy gradient theorem (PGT) [Sutton et al., 2000]

PG method: $\theta^{(k+1)} = \theta^{(k)} - \eta_k \nabla_{\theta} V_{\rho}(\theta^{(k)})$

- State-action cost function (a.k.a Q-function) & advantage function

$$Q_{s,a}(\theta) := \mathbb{E}_{a_t \sim \pi_{s_t}(\theta), s_{t+1} \sim P(\cdot | s_t, a_t)} \left[\sum_{t=0}^{\infty} \gamma^t c(s_t, a_t) \mid s_0 = s, a_0 = a \right]$$

$$A_{s,a}(\theta) := Q_{s,a}(\theta) - \mathbb{E}_{a' \sim \pi_s(\theta)} [Q_{s,a'}(\theta)]$$

$$\begin{aligned} V_{\rho}(\theta) &= \mathbb{E}_{s \sim \rho} [V_s(\theta)] \\ &= \mathbb{E}_{s \sim \rho, a \sim \pi_s(\theta)} [Q_{s,a}(\theta)] \end{aligned}$$

- State visitation distribution of the MDP $d^{\pi}(\rho) \in \Delta(\mathcal{S})$

Policy gradient theorem (PGT) [Sutton et al., 2000]

PG method: $\theta^{(k+1)} = \theta^{(k)} - \eta_k \nabla_{\theta} V_{\rho}(\theta^{(k)})$

- State-action cost function (a.k.a Q-function) & advantage function

$$Q_{s,a}(\theta) := \mathbb{E}_{a_t \sim \pi_{s_t}(\theta), s_{t+1} \sim P(\cdot | s_t, a_t)} \left[\sum_{t=0}^{\infty} \gamma^t c(s_t, a_t) \mid s_0 = s, a_0 = a \right]$$

$$A_{s,a}(\theta) := Q_{s,a}(\theta) - \mathbb{E}_{a' \sim \pi_s(\theta)} [Q_{s,a'}(\theta)]$$

$$\begin{aligned} V_{\rho}(\theta) &= \mathbb{E}_{s \sim \rho} [V_s(\theta)] \\ &= \mathbb{E}_{s \sim \rho, a \sim \pi_s(\theta)} [Q_{s,a}(\theta)] \end{aligned}$$

- State visitation distribution of the MDP $d^{\pi}(\rho) \in \Delta(\mathcal{S})$

$$d_s^{\pi}(\rho) := (1 - \gamma) \mathbb{E}_{s_0 \sim \rho} \left[\sum_{t=0}^{\infty} \gamma^t P(s_t = s \mid s_0, \pi) \right]$$

Policy gradient theorem (PGT) [Sutton et al., 2000]

PG method: $\theta^{(k+1)} = \theta^{(k)} - \eta_k \nabla_{\theta} V_{\rho}(\theta^{(k)})$

- State-action cost function (a.k.a Q-function) & advantage function

$$Q_{s,a}(\theta) := \mathbb{E}_{a_t \sim \pi_{s_t}(\theta), s_{t+1} \sim P(\cdot | s_t, a_t)} \left[\sum_{t=0}^{\infty} \gamma^t c(s_t, a_t) \mid s_0 = s, a_0 = a \right]$$

$$A_{s,a}(\theta) := Q_{s,a}(\theta) - \mathbb{E}_{a' \sim \pi_s(\theta)} [Q_{s,a'}(\theta)]$$

$$\begin{aligned} V_{\rho}(\theta) &= \mathbb{E}_{s \sim \rho} [V_s(\theta)] \\ &= \mathbb{E}_{s \sim \rho, a \sim \pi_s(\theta)} [Q_{s,a}(\theta)] \end{aligned}$$

- State visitation distribution of the MDP $d^{\pi}(\rho) \in \Delta(\mathcal{S})$

$$d_s^{\pi}(\rho) := (1 - \gamma) \mathbb{E}_{s_0 \sim \rho} \left[\sum_{t=0}^{\infty} \gamma^t P(s_t = s \mid s_0, \pi) \right]$$

- PGT

Policy gradient theorem (PGT) [Sutton et al., 2000]

PG method: $\theta^{(k+1)} = \theta^{(k)} - \eta_k \nabla_{\theta} V_{\rho}(\theta^{(k)})$

- State-action cost function (a.k.a Q-function) & advantage function

$$Q_{s,a}(\theta) := \mathbb{E}_{a_t \sim \pi_{s_t}(\theta), s_{t+1} \sim P(\cdot | s_t, a_t)} \left[\sum_{t=0}^{\infty} \gamma^t c(s_t, a_t) \mid s_0 = s, a_0 = a \right]$$

$$A_{s,a}(\theta) := Q_{s,a}(\theta) - \mathbb{E}_{a' \sim \pi_s(\theta)} [Q_{s,a'}(\theta)]$$

$$\begin{aligned} V_{\rho}(\theta) &= \mathbb{E}_{s \sim \rho} [V_s(\theta)] \\ &= \mathbb{E}_{s \sim \rho, a \sim \pi_s(\theta)} [Q_{s,a}(\theta)] \end{aligned}$$

- State visitation distribution of the MDP $d^{\pi}(\rho) \in \Delta(\mathcal{S})$

$$d_s^{\pi}(\rho) := (1 - \gamma) \mathbb{E}_{s_0 \sim \rho} \left[\sum_{t=0}^{\infty} \gamma^t P(s_t = s \mid s_0, \pi) \right]$$

- PGT

$$\nabla_{\theta} V_{\rho}(\theta) = \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d^{\pi}(\rho), a \sim \pi_s(\theta)} [Q_{s,a}(\theta) \nabla_{\theta} \log \pi_{s,a}(\theta)]$$

Derivation of Policy Gradient Theorem

$$\nabla_{\theta} V_{\rho}(\theta) = \frac{1}{1-\gamma} \mathbb{E}_{s \sim d^{\pi(\theta)}(\rho), a \sim \pi_s(\theta)} [Q_{s,a}(\theta) \nabla_{\theta} \log \pi_{s,a}(\theta)]$$

Derivation of Policy Gradient Theorem

$$\nabla_{\theta} V_{\rho}(\theta) = \frac{1}{1-\gamma} \mathbb{E}_{s \sim d^{\pi(\theta)}(\rho), a \sim \pi_s(\theta)} [Q_{s,a}(\theta) \nabla_{\theta} \log \pi_{s,a}(\theta)]$$

- *Proof:*

Derivation of Policy Gradient Theorem

$$\nabla_{\theta} V_{\rho}(\theta) = \frac{1}{1-\gamma} \mathbb{E}_{s \sim d^{\pi(\theta)}(\rho), a \sim \pi_s(\theta)} [Q_{s,a}(\theta) \nabla_{\theta} \log \pi_{s,a}(\theta)]$$

• *Proof:*
$$\nabla_{\theta} V_{\rho}(\theta) = \nabla_{\theta} \sum_{s_0 \in \mathcal{S}, a_0 \in \mathcal{A}} \rho(s_0) \pi_{s_0, a_0}(\theta) Q_{s_0, a_0}(\theta)$$

Derivation of Policy Gradient Theorem

$$\nabla_{\theta} V_{\rho}(\theta) = \frac{1}{1-\gamma} \mathbb{E}_{s \sim d^{\pi(\theta)}(\rho), a \sim \pi_s(\theta)} [Q_{s,a}(\theta) \nabla_{\theta} \log \pi_{s,a}(\theta)]$$

• *Proof:*

$$\begin{aligned} \nabla_{\theta} V_{\rho}(\theta) &= \nabla_{\theta} \sum_{s_0 \in \mathcal{S}, a_0 \in \mathcal{A}} \rho(s_0) \pi_{s_0, a_0}(\theta) Q_{s_0, a_0}(\theta) \\ &= \sum_{s_0, a_0} \rho(s_0) (\nabla_{\theta} \pi_{s_0, a_0}(\theta)) Q_{s_0, a_0}(\theta) + \sum_{s_0, a_0} \rho(s_0) \pi_{s_0, a_0}(\theta) \nabla_{\theta} Q_{s_0, a_0}(\theta) \end{aligned}$$

Derivation of Policy Gradient Theorem

$$\nabla_{\theta} V_{\rho}(\theta) = \frac{1}{1-\gamma} \mathbb{E}_{s \sim d^{\pi(\theta)}(\rho), a \sim \pi_s(\theta)} [Q_{s,a}(\theta) \nabla_{\theta} \log \pi_{s,a}(\theta)]$$

• *Proof:*

$$\begin{aligned} \nabla_{\theta} V_{\rho}(\theta) &= \nabla_{\theta} \sum_{s_0 \in \mathcal{S}, a_0 \in \mathcal{A}} \rho(s_0) \pi_{s_0, a_0}(\theta) Q_{s_0, a_0}(\theta) \\ &= \sum_{s_0, a_0} \rho(s_0) (\nabla_{\theta} \pi_{s_0, a_0}(\theta)) Q_{s_0, a_0}(\theta) + \sum_{s_0, a_0} \rho(s_0) \pi_{s_0, a_0}(\theta) \nabla_{\theta} Q_{s_0, a_0}(\theta) \\ &= \sum_{s_0, a_0} \rho(s_0) \pi_{s_0, a_0}(\theta) (\nabla_{\theta} \log \pi_{s_0, a_0}(\theta)) Q_{s_0, a_0}(\theta) \\ &\quad + \sum_{s_0, a_0} \rho(s_0) \pi_{s_0, a_0}(\theta) \nabla_{\theta} (c(s_0, a_0) + \gamma \sum_{s_1} P(s_1 | s_0, a_0) V_{s_1}(\theta)) \end{aligned}$$

Derivation of Policy Gradient Theorem

$$\nabla_{\theta} V_{\rho}(\theta) = \frac{1}{1-\gamma} \mathbb{E}_{s \sim d^{\pi(\theta)}(\rho), a \sim \pi_s(\theta)} [Q_{s,a}(\theta) \nabla_{\theta} \log \pi_{s,a}(\theta)]$$

• *Proof:*

$$\begin{aligned} \nabla_{\theta} V_{\rho}(\theta) &= \nabla_{\theta} \sum_{s_0 \in \mathcal{S}, a_0 \in \mathcal{A}} \rho(s_0) \pi_{s_0, a_0}(\theta) Q_{s_0, a_0}(\theta) \\ &= \sum_{s_0, a_0} \rho(s_0) (\nabla_{\theta} \pi_{s_0, a_0}(\theta)) Q_{s_0, a_0}(\theta) + \sum_{s_0, a_0} \rho(s_0) \pi_{s_0, a_0}(\theta) \nabla_{\theta} Q_{s_0, a_0}(\theta) \\ &= \sum_{s_0, a_0} \rho(s_0) \pi_{s_0, a_0}(\theta) (\nabla_{\theta} \log \pi_{s_0, a_0}(\theta)) Q_{s_0, a_0}(\theta) \\ &\quad + \sum_{s_0, a_0} \rho(s_0) \pi_{s_0, a_0}(\theta) \nabla_{\theta} \left(c(s_0, a_0) + \gamma \sum_{s_1} P(s_1 | s_0, a_0) V_{s_1}(\theta) \right) \end{aligned}$$

Bellman Equation

Derivation of Policy Gradient Theorem

$$\nabla_{\theta} V_{\rho}(\theta) = \frac{1}{1-\gamma} \mathbb{E}_{s \sim d^{\pi(\theta)}(\rho), a \sim \pi_s(\theta)} [Q_{s,a}(\theta) \nabla_{\theta} \log \pi_{s,a}(\theta)]$$

• *Proof:*

$$\begin{aligned} \nabla_{\theta} V_{\rho}(\theta) &= \nabla_{\theta} \sum_{s_0 \in \mathcal{S}, a_0 \in \mathcal{A}} \rho(s_0) \pi_{s_0, a_0}(\theta) Q_{s_0, a_0}(\theta) \\ &= \sum_{s_0, a_0} \rho(s_0) \left(\nabla_{\theta} \pi_{s_0, a_0}(\theta) \right) Q_{s_0, a_0}(\theta) + \sum_{s_0, a_0} \rho(s_0) \pi_{s_0, a_0}(\theta) \nabla_{\theta} Q_{s_0, a_0}(\theta) \\ &= \sum_{s_0, a_0} \rho(s_0) \pi_{s_0, a_0}(\theta) \left(\nabla_{\theta} \log \pi_{s_0, a_0}(\theta) \right) Q_{s_0, a_0}(\theta) \\ &\quad + \sum_{s_0, a_0} \rho(s_0) \pi_{s_0, a_0}(\theta) \nabla_{\theta} \left(c(s_0, a_0) + \gamma \sum_{s_1} P(s_1 | s_0, a_0) V_{s_1}(\theta) \right) \\ &= \sum_{s_0, a_0} \rho(s_0) \pi_{s_0, a_0}(\theta) \left(\nabla_{\theta} \log \pi_{s_0, a_0}(\theta) \right) Q_{s_0, a_0}(\theta) \\ &\quad + \gamma \sum_{s_0, a_0, s_1} \rho(s_0) \pi_{s_0, a_0}(\theta) P(s_1 | s_0, a_0) \nabla_{\theta} V_{s_1}(\theta) \end{aligned}$$

Bellman Equation

Derivation of Policy Gradient Theorem

$$\nabla_{\theta} V_{\rho}(\theta) = \frac{1}{1-\gamma} \mathbb{E}_{s \sim d^{\pi(\theta)}(\rho), a \sim \pi_s(\theta)} [Q_{s,a}(\theta) \nabla_{\theta} \log \pi_{s,a}(\theta)]$$

• *Proof:*

$$\begin{aligned} \nabla_{\theta} V_{\rho}(\theta) &= \nabla_{\theta} \sum_{s_0 \in \mathcal{S}, a_0 \in \mathcal{A}} \rho(s_0) \pi_{s_0, a_0}(\theta) Q_{s_0, a_0}(\theta) \\ &= \sum_{s_0, a_0} \rho(s_0) \left(\nabla_{\theta} \pi_{s_0, a_0}(\theta) \right) Q_{s_0, a_0}(\theta) + \sum_{s_0, a_0} \rho(s_0) \pi_{s_0, a_0}(\theta) \nabla_{\theta} Q_{s_0, a_0}(\theta) \\ &= \sum_{s_0, a_0} \rho(s_0) \pi_{s_0, a_0}(\theta) \left(\nabla_{\theta} \log \pi_{s_0, a_0}(\theta) \right) Q_{s_0, a_0}(\theta) \\ &\quad + \sum_{s_0, a_0} \rho(s_0) \pi_{s_0, a_0}(\theta) \nabla_{\theta} \left(c(s_0, a_0) + \gamma \sum_{s_1} P(s_1 | s_0, a_0) V_{s_1}(\theta) \right) \quad \text{Bellman Equation} \\ &= \sum_{s_0, a_0} \rho(s_0) \pi_{s_0, a_0}(\theta) \left(\nabla_{\theta} \log \pi_{s_0, a_0}(\theta) \right) Q_{s_0, a_0}(\theta) \\ &\quad + \gamma \sum_{s_0, a_0, s_1} \rho(s_0) \pi_{s_0, a_0}(\theta) P(s_1 | s_0, a_0) \nabla_{\theta} V_{s_1}(\theta) \\ &= \mathbb{E} \left[Q_{s_0, a_0}(\theta) \nabla_{\theta} \log \pi_{s_0, a_0}(\theta) \right] + \gamma \mathbb{E} \left[\nabla_{\theta} V_{s_1}(\theta) \right] \end{aligned}$$

Derivation of Policy Gradient Theorem

$$\nabla_{\theta} V_{\rho}(\theta) = \frac{1}{1-\gamma} \mathbb{E}_{s \sim d^{\pi(\theta)}(\rho), a \sim \pi_s(\theta)} [Q_{s,a}(\theta) \nabla_{\theta} \log \pi_{s,a}(\theta)]$$

• *Proof:*

$$\begin{aligned} \nabla_{\theta} V_{\rho}(\theta) &= \nabla_{\theta} \sum_{s_0 \in \mathcal{S}, a_0 \in \mathcal{A}} \rho(s_0) \pi_{s_0, a_0}(\theta) Q_{s_0, a_0}(\theta) \\ &= \sum_{s_0, a_0} \rho(s_0) \left(\nabla_{\theta} \pi_{s_0, a_0}(\theta) \right) Q_{s_0, a_0}(\theta) + \sum_{s_0, a_0} \rho(s_0) \pi_{s_0, a_0}(\theta) \nabla_{\theta} Q_{s_0, a_0}(\theta) \\ &= \sum_{s_0, a_0} \rho(s_0) \pi_{s_0, a_0}(\theta) \left(\nabla_{\theta} \log \pi_{s_0, a_0}(\theta) \right) Q_{s_0, a_0}(\theta) \\ &\quad + \sum_{s_0, a_0} \rho(s_0) \pi_{s_0, a_0}(\theta) \nabla_{\theta} \left(c(s_0, a_0) + \gamma \sum_{s_1} P(s_1 | s_0, a_0) V_{s_1}(\theta) \right) \quad \text{Bellman Equation} \\ &= \sum_{s_0, a_0} \rho(s_0) \pi_{s_0, a_0}(\theta) \left(\nabla_{\theta} \log \pi_{s_0, a_0}(\theta) \right) Q_{s_0, a_0}(\theta) \\ &\quad + \gamma \sum_{s_0, a_0, s_1} \rho(s_0) \pi_{s_0, a_0}(\theta) P(s_1 | s_0, a_0) \nabla_{\theta} V_{s_1}(\theta) \\ &= \mathbb{E} \left[Q_{s_0, a_0}(\theta) \nabla_{\theta} \log \pi_{s_0, a_0}(\theta) \right] + \gamma \mathbb{E} \left[\nabla_{\theta} V_{s_1}(\theta) \right] \\ &= \mathbb{E} \left[Q_{s_0, a_0}(\theta) \nabla_{\theta} \log \pi_{s_0, a_0}(\theta) \right] + \gamma \mathbb{E} \left[Q_{s_1, a_1}(\theta) \nabla_{\theta} \log \pi_{s_1, a_1}(\theta) \right] + \dots \end{aligned}$$

Derivation of Policy Gradient Theorem

$$\nabla_{\theta} V_{\rho}(\theta) = \frac{1}{1-\gamma} \mathbb{E}_{s \sim d^{\pi(\theta)}(\rho), a \sim \pi_s(\theta)} [Q_{s,a}(\theta) \nabla_{\theta} \log \pi_{s,a}(\theta)]$$

• *Proof:*

$$\begin{aligned} \nabla_{\theta} V_{\rho}(\theta) &= \nabla_{\theta} \sum_{s_0 \in \mathcal{S}, a_0 \in \mathcal{A}} \rho(s_0) \pi_{s_0, a_0}(\theta) Q_{s_0, a_0}(\theta) \\ &= \sum_{s_0, a_0} \rho(s_0) \left(\nabla_{\theta} \pi_{s_0, a_0}(\theta) \right) Q_{s_0, a_0}(\theta) + \sum_{s_0, a_0} \rho(s_0) \pi_{s_0, a_0}(\theta) \nabla_{\theta} Q_{s_0, a_0}(\theta) \\ &= \sum_{s_0, a_0} \rho(s_0) \pi_{s_0, a_0}(\theta) \left(\nabla_{\theta} \log \pi_{s_0, a_0}(\theta) \right) Q_{s_0, a_0}(\theta) \\ &\quad + \sum_{s_0, a_0} \rho(s_0) \pi_{s_0, a_0}(\theta) \nabla_{\theta} \left(c(s_0, a_0) + \gamma \sum_{s_1} P(s_1 | s_0, a_0) V_{s_1}(\theta) \right) \quad \text{Bellman Equation} \\ &= \sum_{s_0, a_0} \rho(s_0) \pi_{s_0, a_0}(\theta) \left(\nabla_{\theta} \log \pi_{s_0, a_0}(\theta) \right) Q_{s_0, a_0}(\theta) \\ &\quad + \gamma \sum_{s_0, a_0, s_1} \rho(s_0) \pi_{s_0, a_0}(\theta) P(s_1 | s_0, a_0) \nabla_{\theta} V_{s_1}(\theta) \\ &= \mathbb{E} \left[Q_{s_0, a_0}(\theta) \nabla_{\theta} \log \pi_{s_0, a_0}(\theta) \right] + \gamma \mathbb{E} \left[\nabla_{\theta} V_{s_1}(\theta) \right] \\ &= \mathbb{E} \left[Q_{s_0, a_0}(\theta) \nabla_{\theta} \log \pi_{s_0, a_0}(\theta) \right] + \gamma \mathbb{E} \left[Q_{s_1, a_1}(\theta) \nabla_{\theta} \log \pi_{s_1, a_1}(\theta) \right] + \dots \\ &= \frac{1}{1-\gamma} \mathbb{E}_{s \sim d^{\pi(\theta)}(\rho), a \sim \pi_s(\theta)} [Q_{s,a}(\theta) \nabla_{\theta} \log \pi_{s,a}(\theta)] \end{aligned}$$