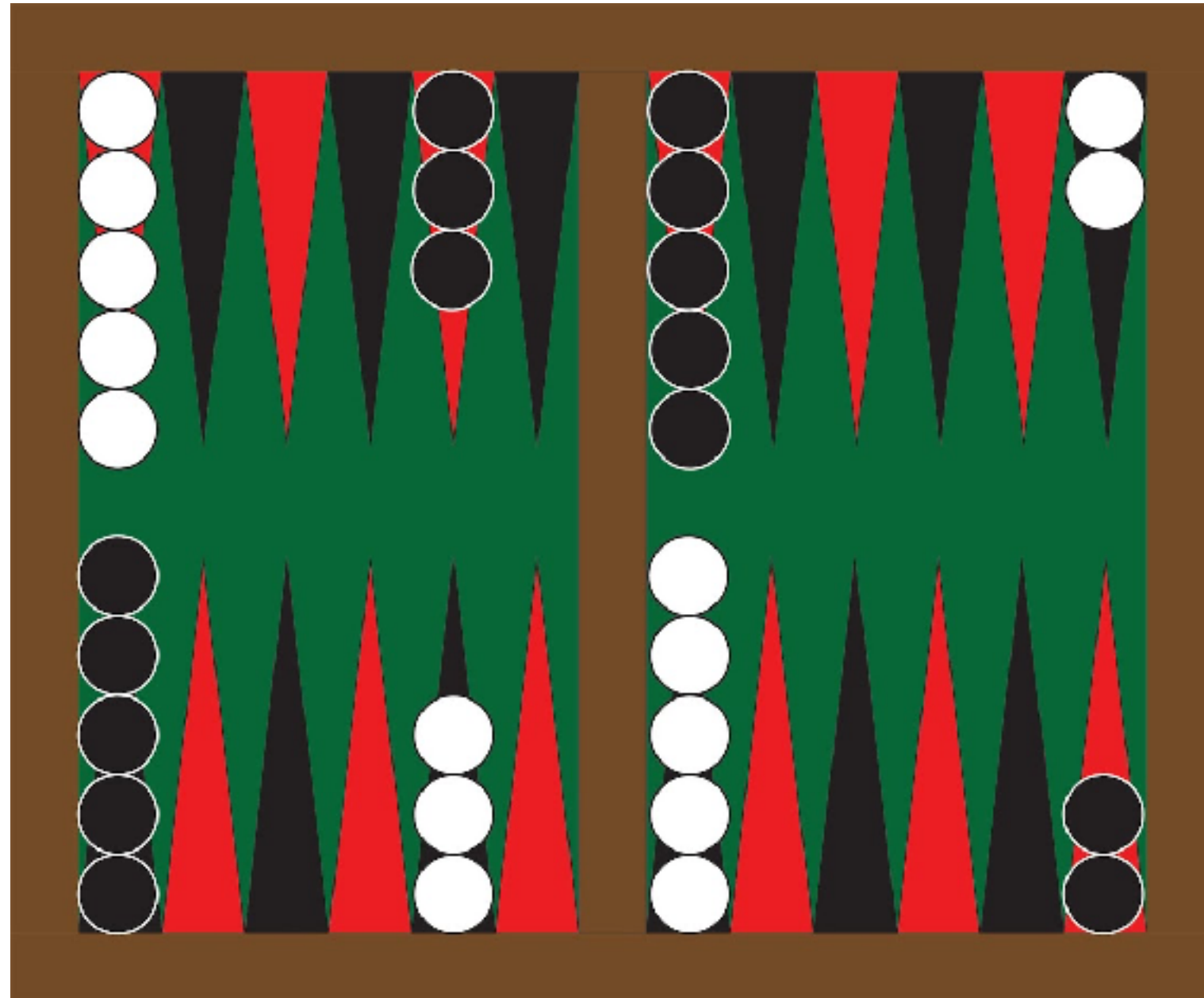


Introduction and Basics of Markov Decision Process

Sham Kakade and Wen Sun

CS 6789: Foundations of Reinforcement Learning

Progress of RL in Practice



TD GAMMON [Tesauro 95]



[AlphaZero, Silver et.al, 17]

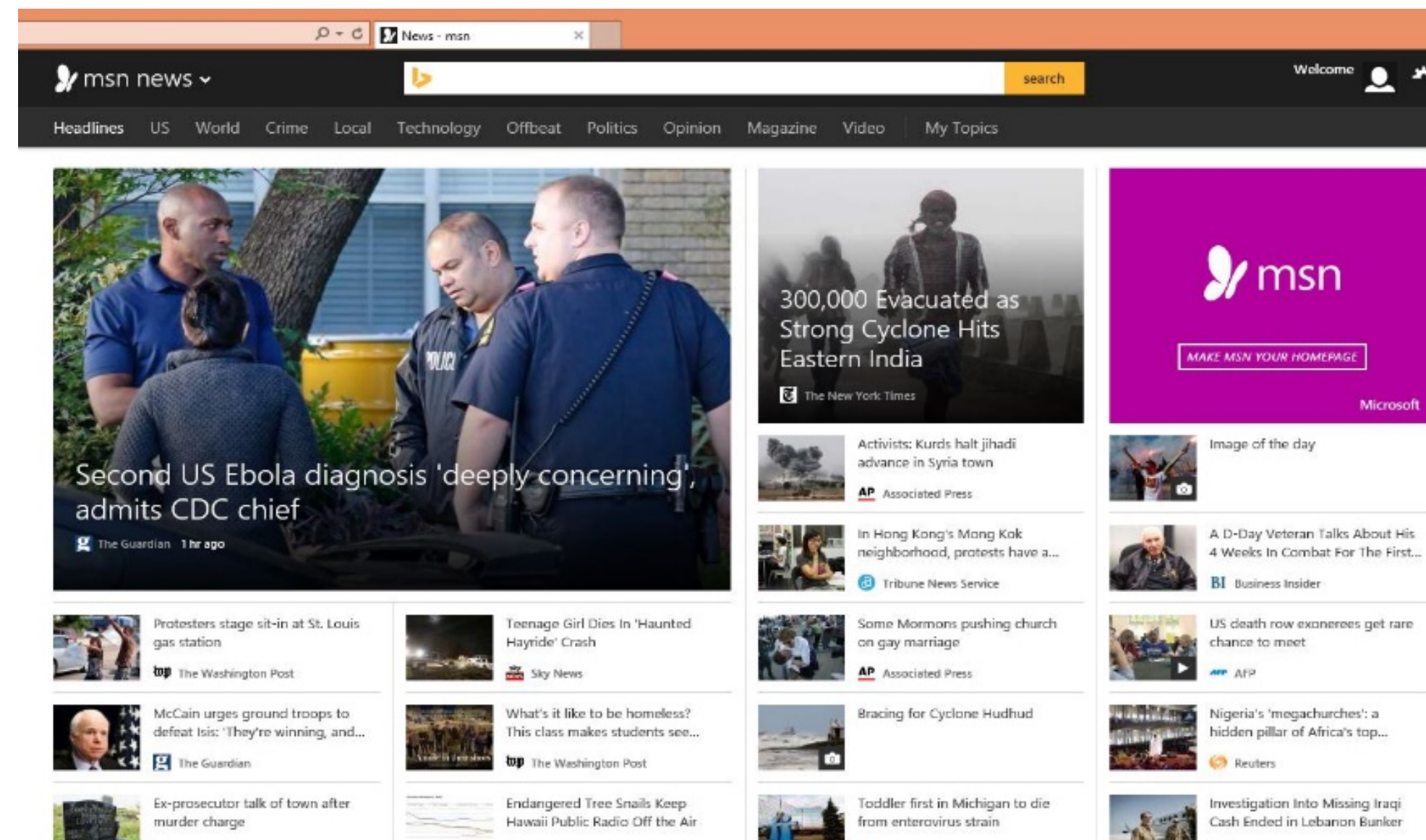


[OpenAI Five, 18]

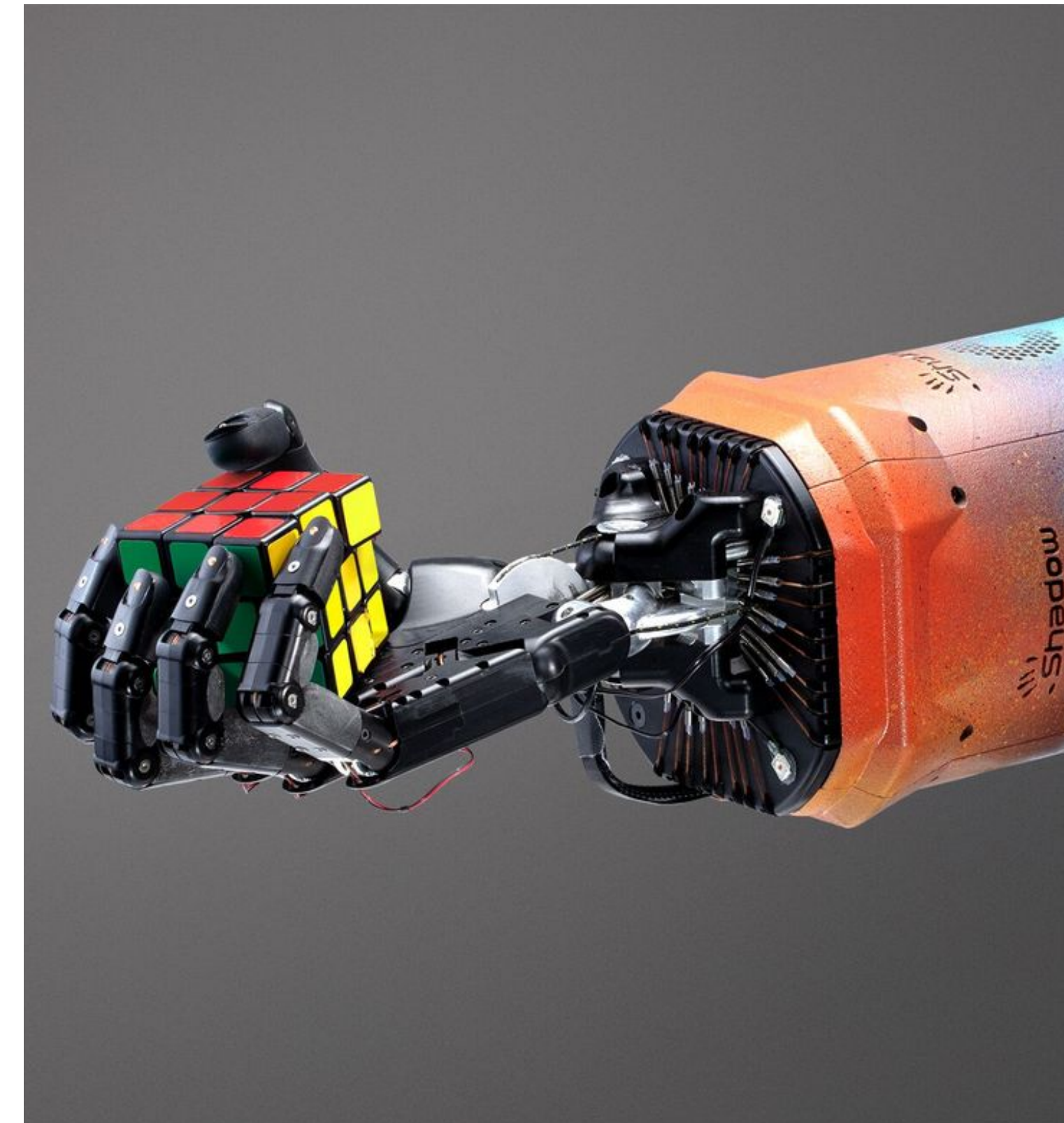
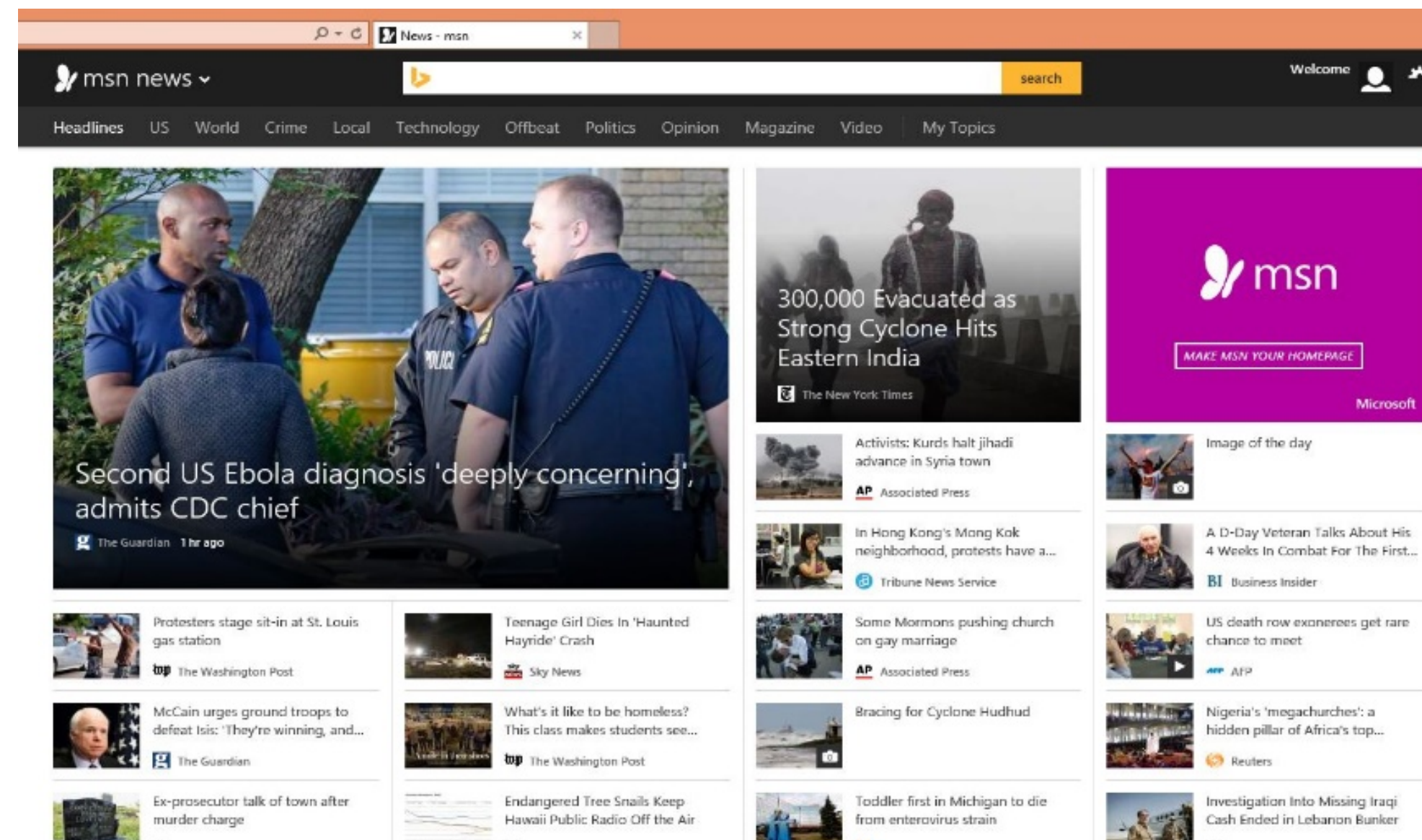
RL in Real World:



RL in Real World:



RL in Real World:



This course focuses on RL Theory

When and Why RL works!

(Convergence, sample / computation complexity, etc)

Four main themes we will cover in this course:

1. Fundamentals (MDPs, statistical limit, lower bounds)
2. Exploration (sample complexity)
3. Policy Gradient (global convergence)
4. Control & Imitation Learning (i.e., learning from demonstrations)

Logistics

Four (HW0-HW3) assignments (total 55%), Course Project (40%), Reading (5%)

(HW0 10%, HW1-3 15% each)

HW0 is out today and due in two week

Prerequisites (HW0)

Deep understanding of Machine Learning, Optimization, Statistics

ML: sample complexity analysis for supervised learning (PAC)

Opt: Convex (linear) optimization, e.g., gradient descent for convex functions

Stats: basics of concentration (e.g., Hoeffding's), tricks such as union bound

Prerequisites (HW0)

Deep understanding of Machine Learning, Optimization, Statistics

ML: sample complexity analysis for supervised learning (PAC)

Opt: Convex (linear) optimization, e.g., gradient descent for convex functions

Stats: basics of concentration (e.g., Hoeffding's), tricks such as union bound

Check out HW0 asap!

Course projects (40%)

- Team work: size 3
- Midterm report (5%), Final presentation (15%), and Final report (20%)
- Basics: **survey** of a set of similar RL theory papers. Reproduce analysis and provide a coherent story
- Advanced: **identify** extensions of existing RL papers, **formulate** theory questions, and **provide** proofs

Course Notes:

Reinforcement Learning Theory & Algorithms

- Book website: <https://rltheorybook.github.io/>
- Many lectures will correspond to chapters in Version 2.
- Reading assignment (5%) is from this book
- Please let us know if you find typos/errors in the book!
We appreciate it!

Outline

1. Definition of infinite horizon discounted MDPs

2. Bellman Optimality

3. State-action distribution

Supervised Learning

Supervised Learning

Given i.i.d examples at training:



(,cat)



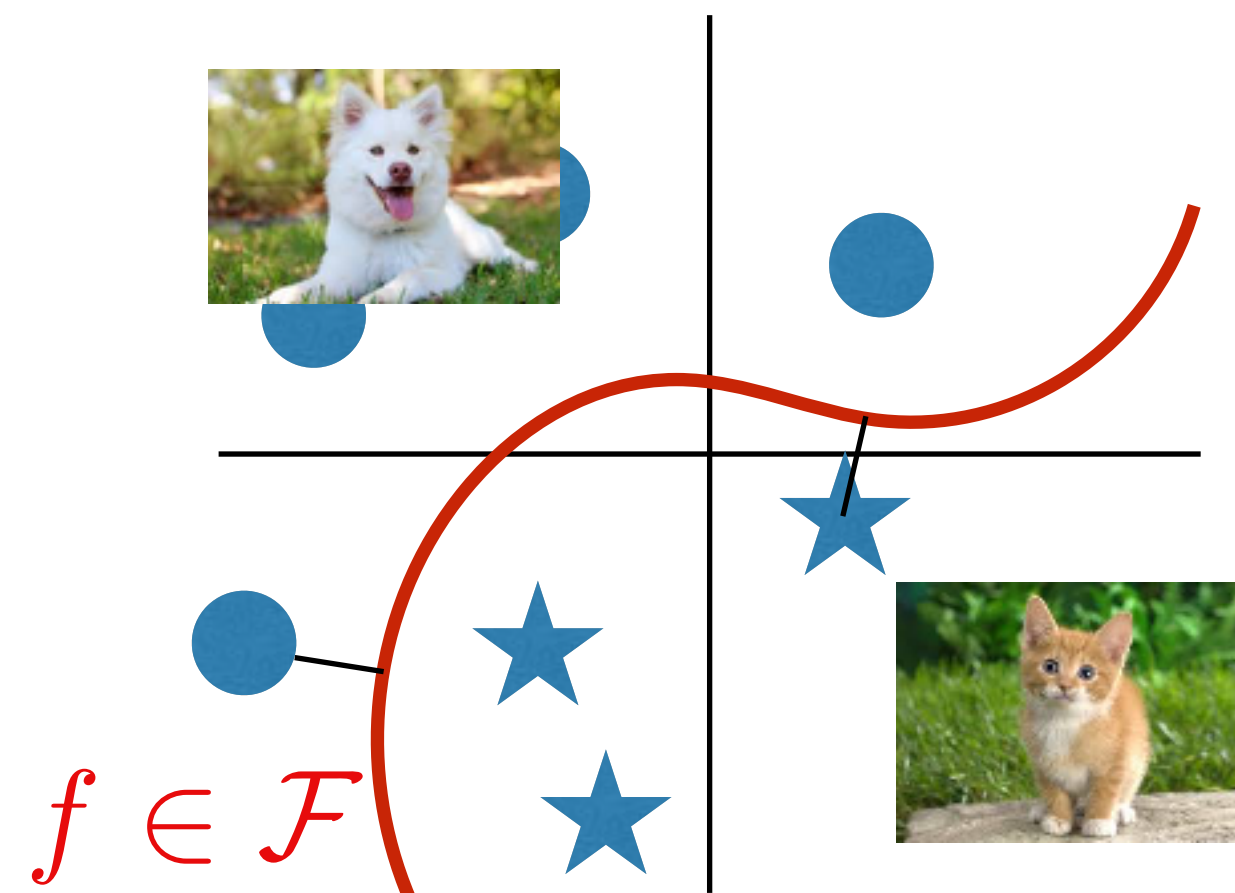
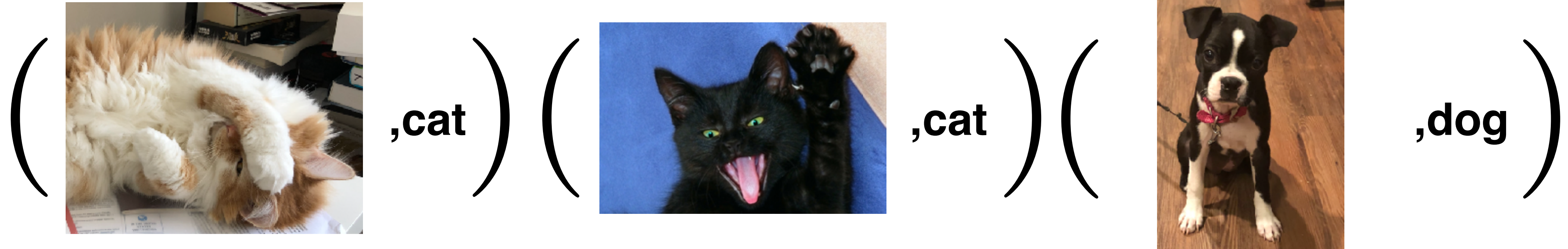
(,cat)



(,dog)

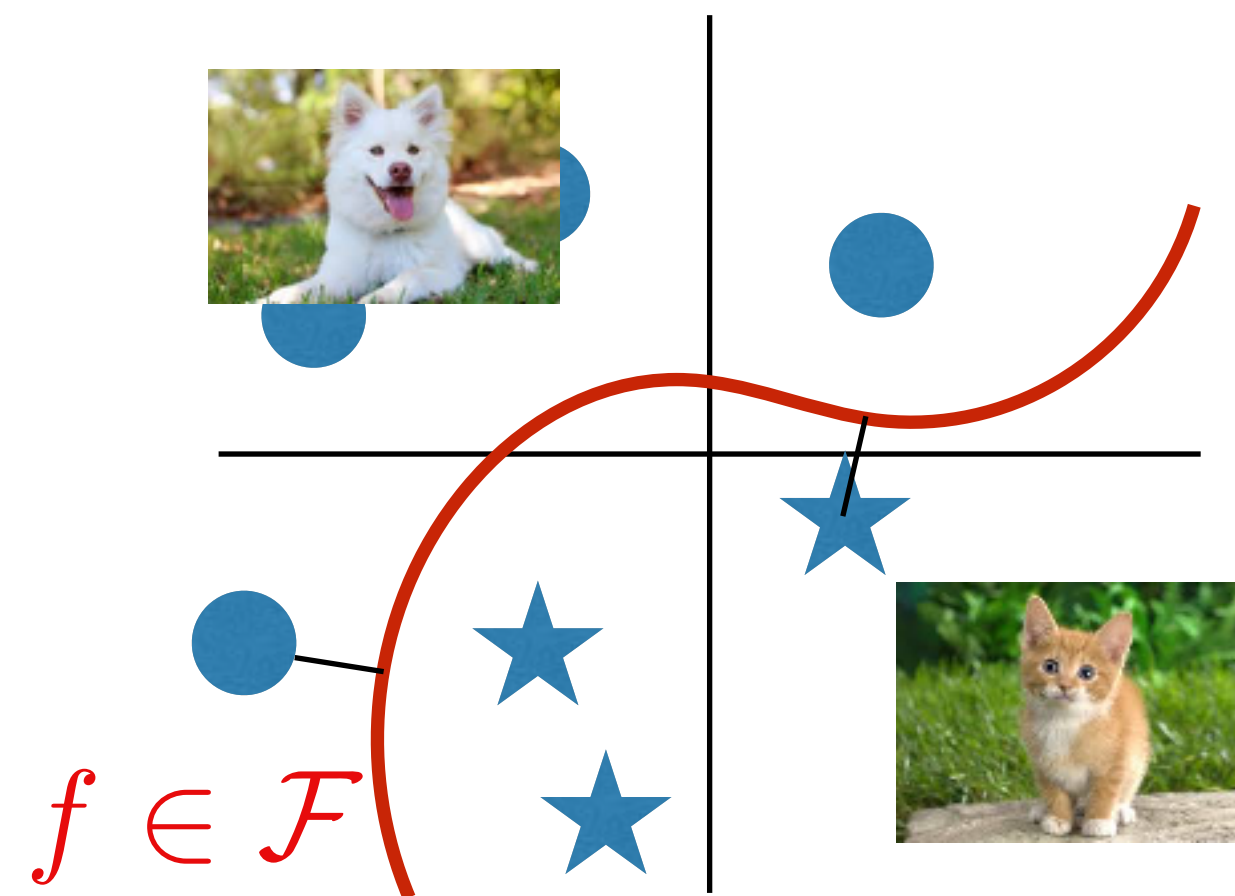
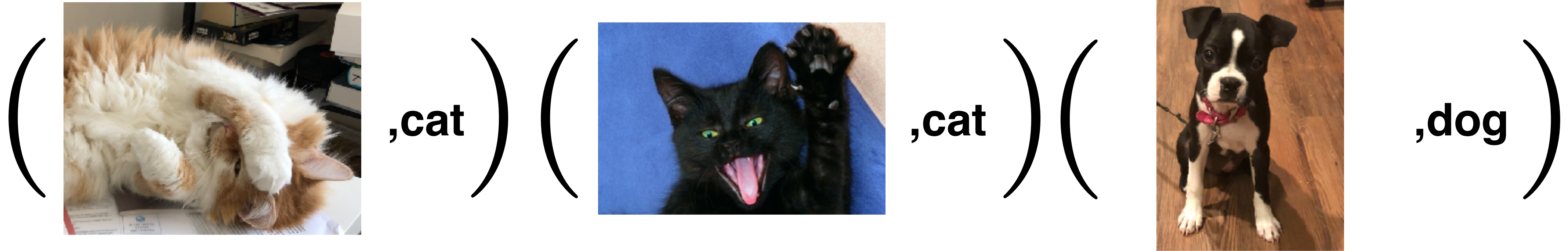
Supervised Learning

Given i.i.d examples at training:

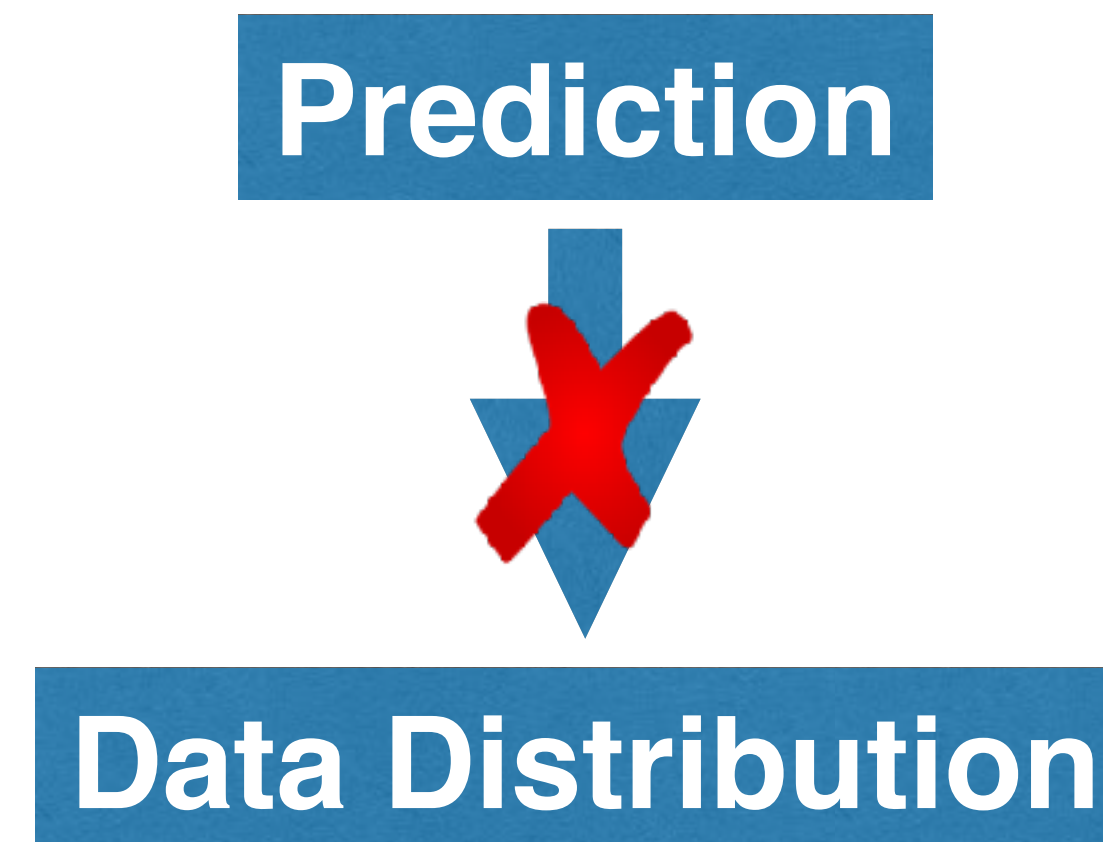


Supervised Learning

Given i.i.d examples at training:

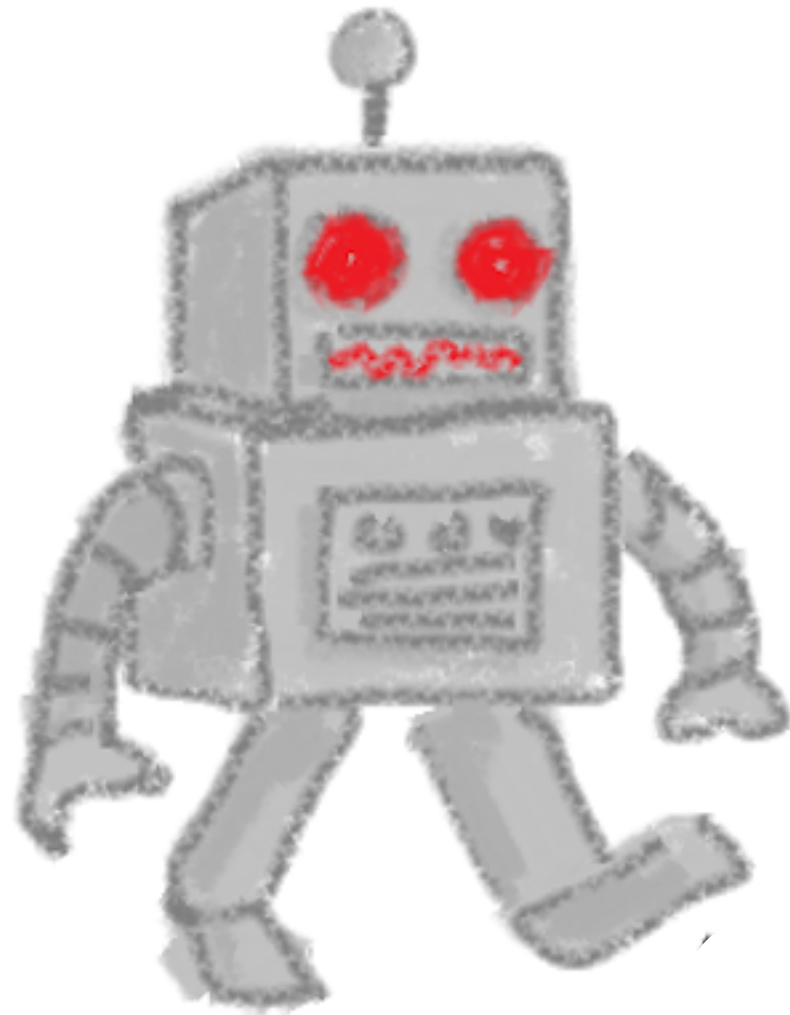


Passive:



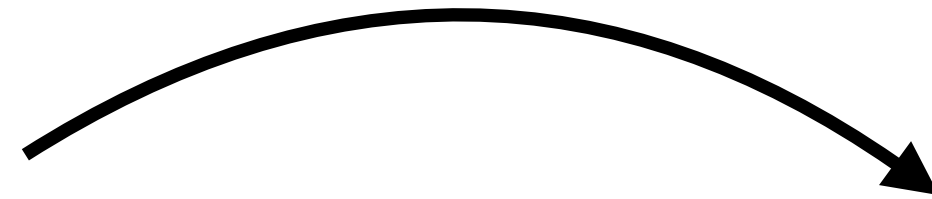
Markov Decision Process

Learning Agent

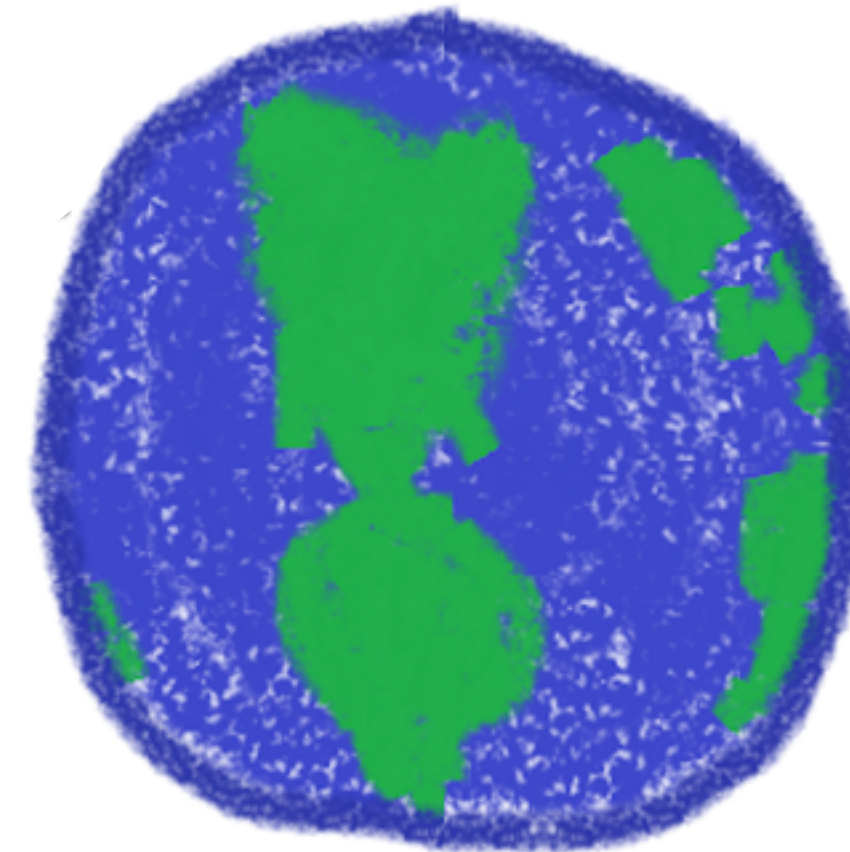


$$a \sim \pi(s)$$

Policy: determine **action** based on **state**

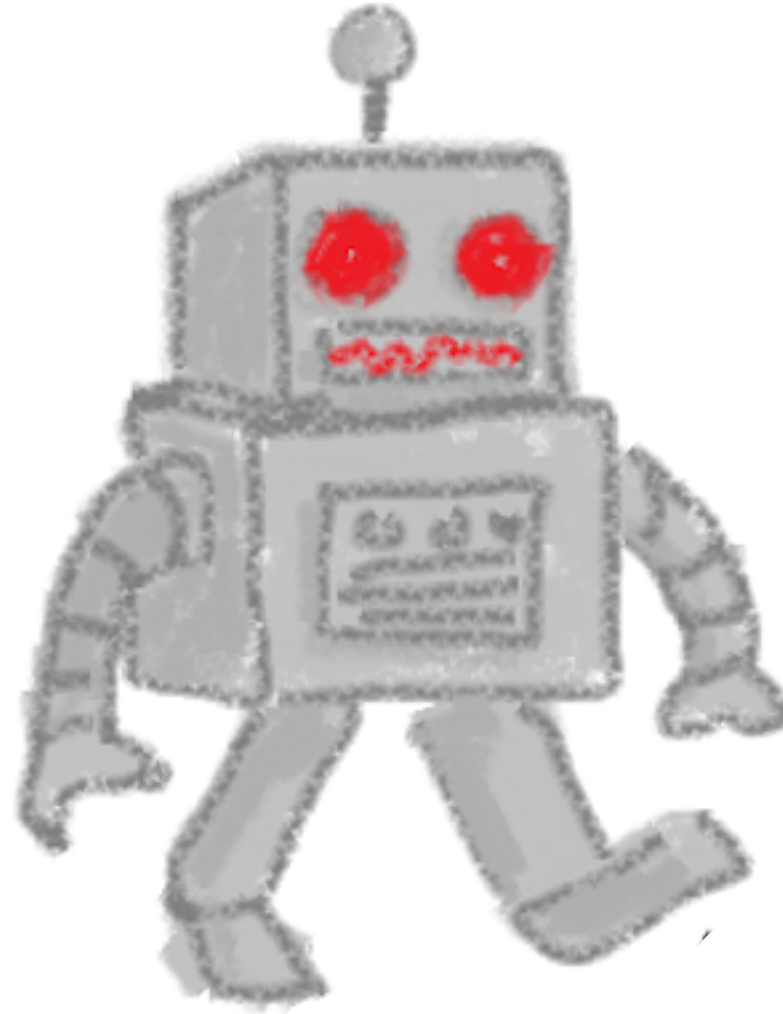


Environment



Markov Decision Process

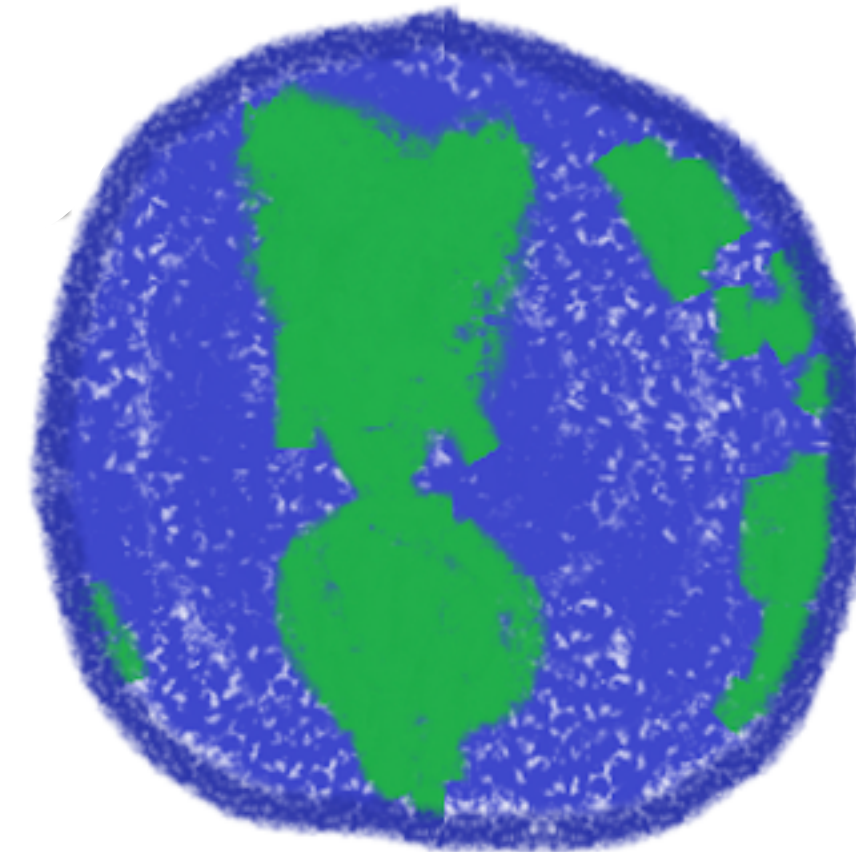
Learning Agent



$$a \sim \pi(s)$$

Policy: determine **action** based on **state**

Environment

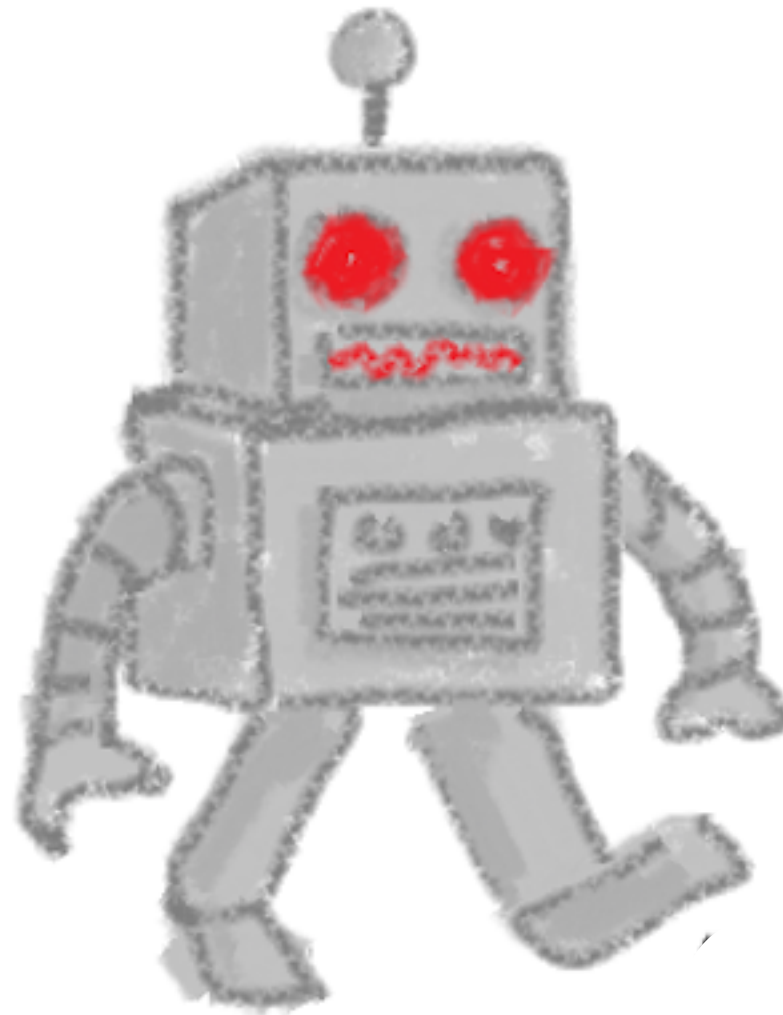


Send **reward** and **next state** from a Markovian transition dynamics

$$r(s, a), s' \sim P(\cdot | s, a)$$

Markov Decision Process

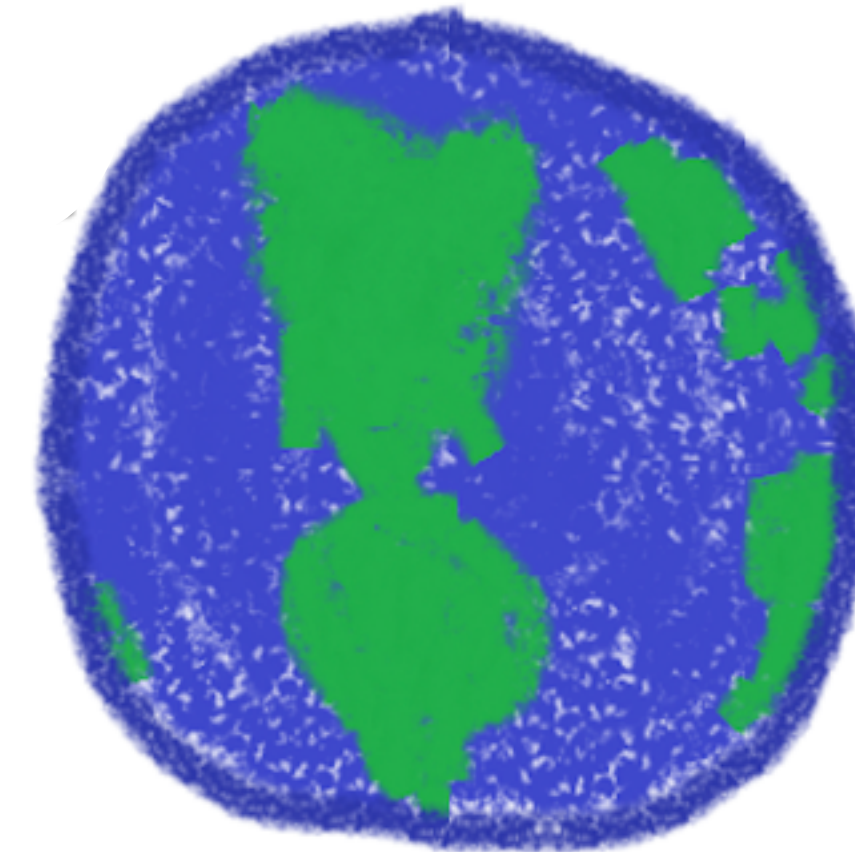
Learning Agent



$$a \sim \pi(s)$$

Policy: determine **action** based on **state**

Environment



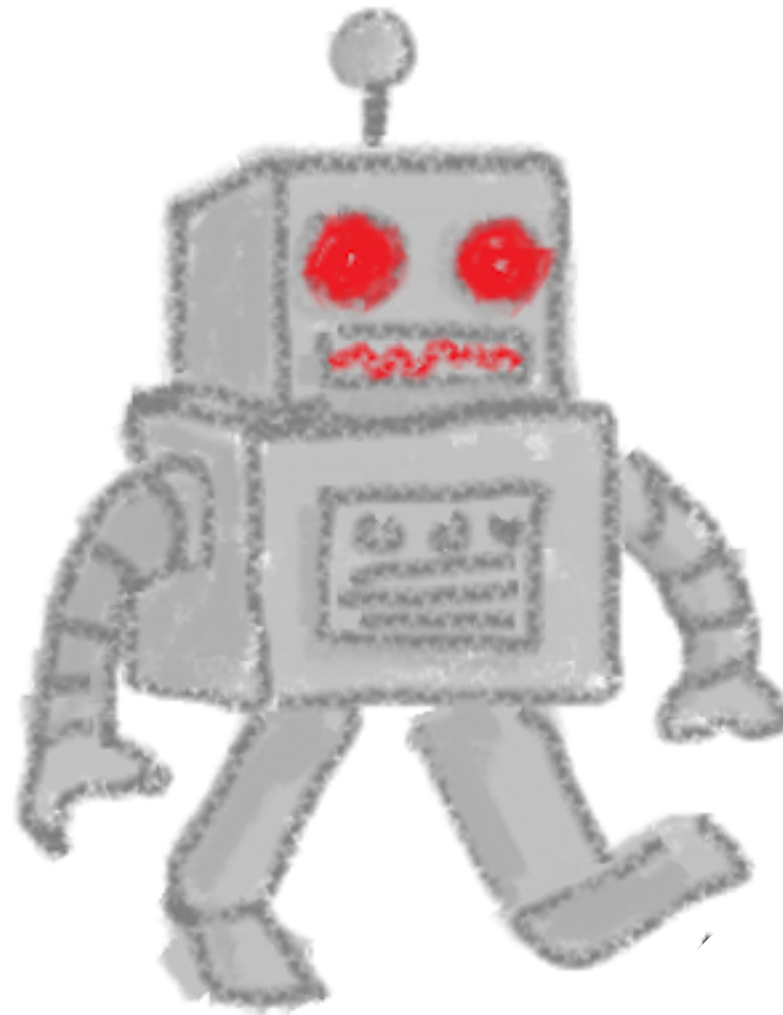
Multiple Steps

Send **reward** and **next state** from a Markovian transition dynamics

$$r(s, a), s' \sim P(\cdot | s, a)$$

Markov Decision Process

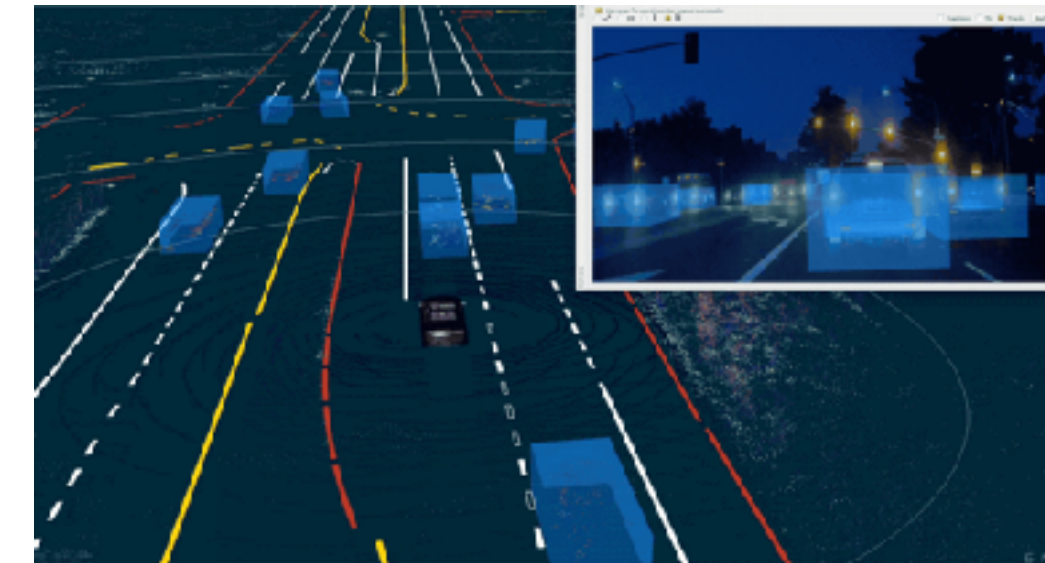
Learning Agent



$$a \sim \pi(s)$$

Policy: determine **action** based on **state**

Environment



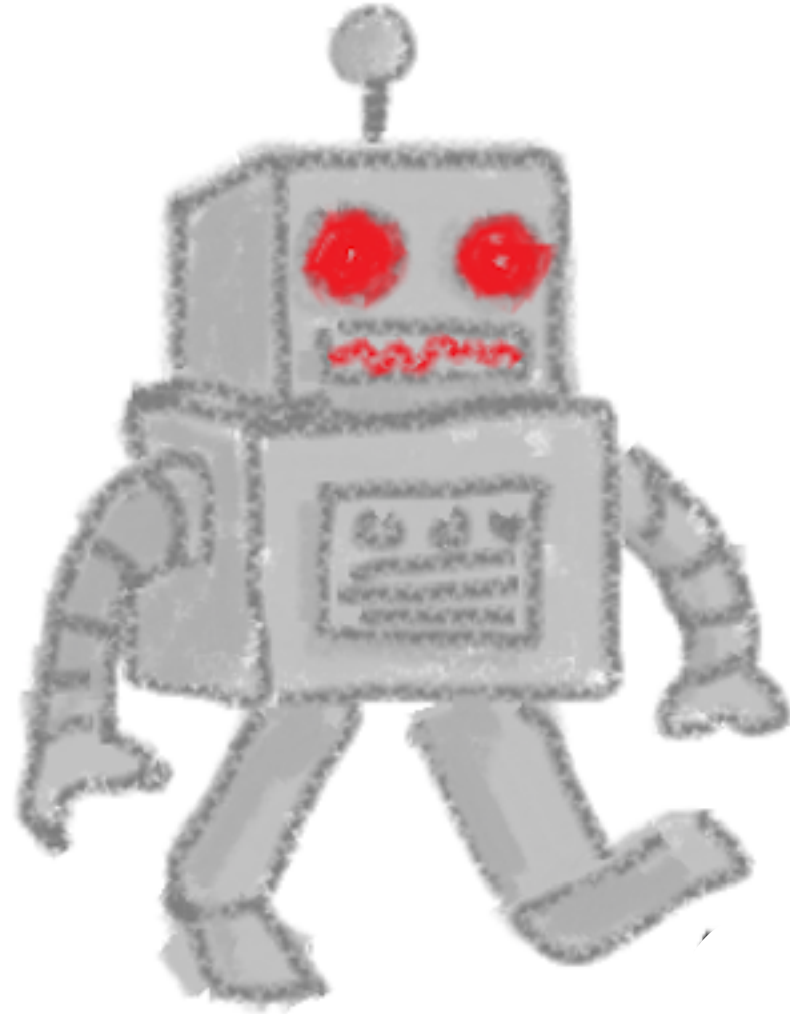
Multiple Steps

Send **reward** and **next state** from a Markovian transition dynamics

$$r(s, a), s' \sim P(\cdot | s, a)$$

Markov Decision Process

Learning Agent



$$a \sim \pi(s)$$

Policy: determine **action** based on **state**

Multiple Steps

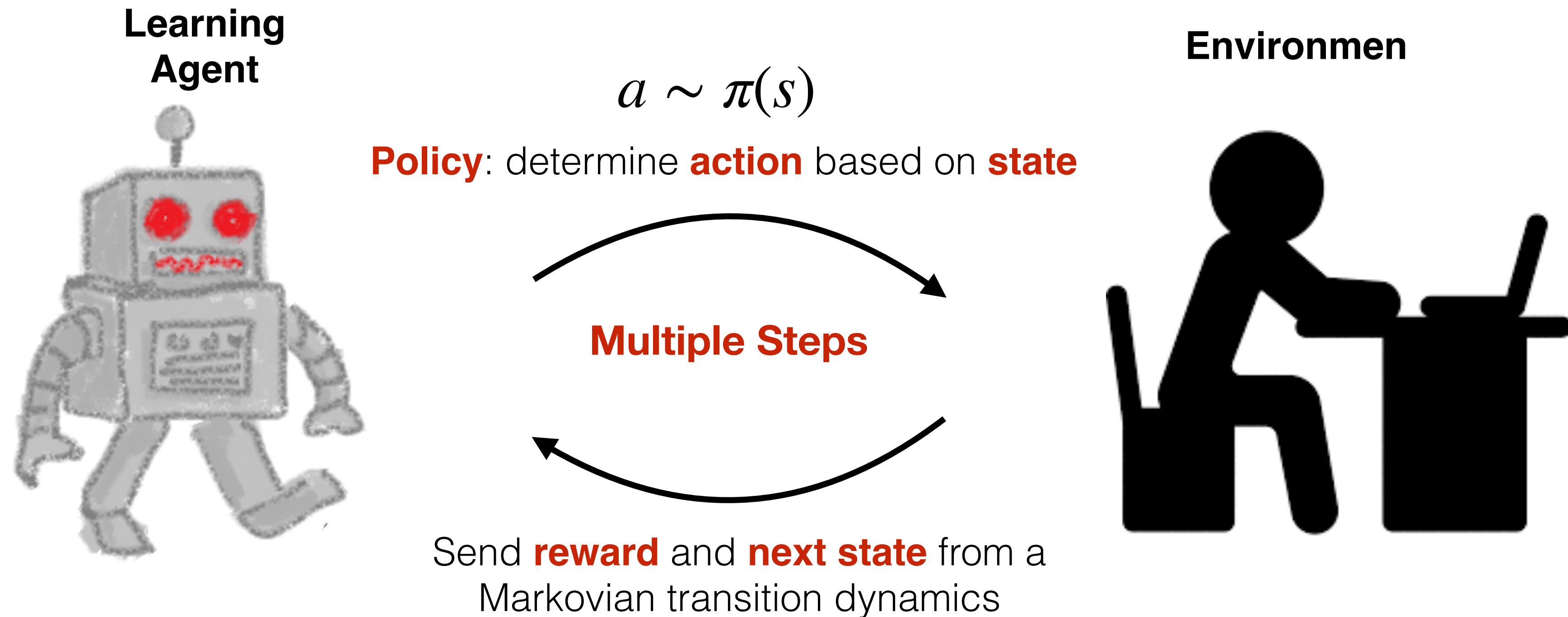
Send **reward** and **next state** from a Markovian transition dynamics

$$r(s, a), s' \sim P(\cdot | s, a)$$

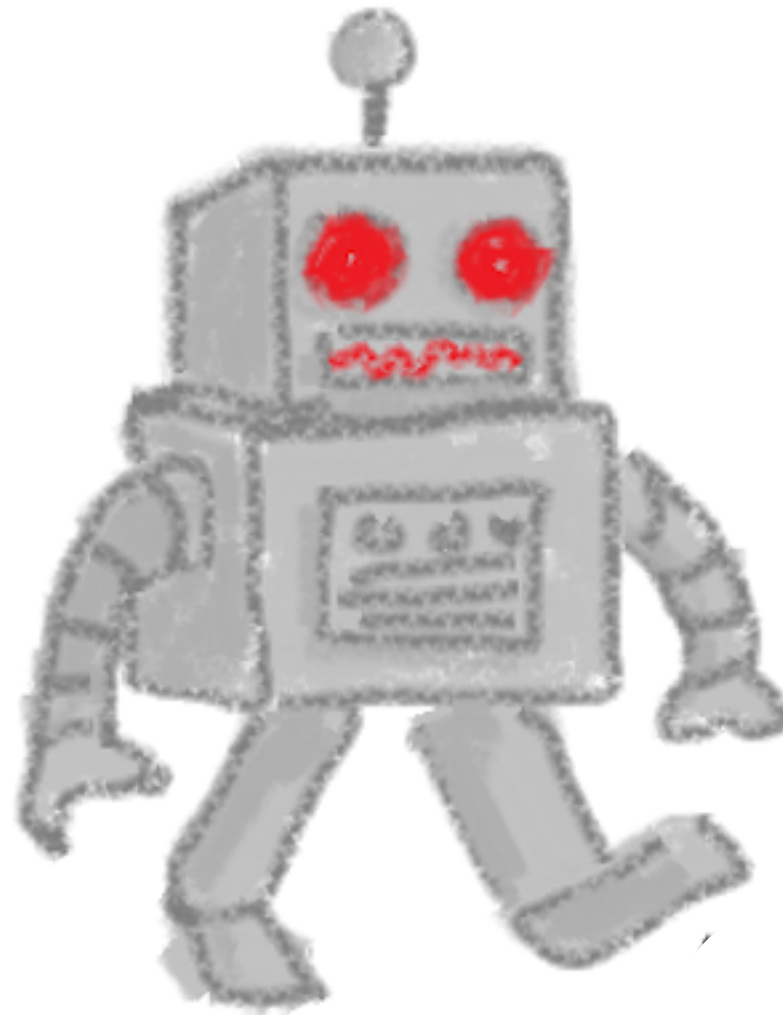
Environment



Markov Decision Process



Learning Agent



$$a \sim \pi(s)$$

Policy: determine **action** based on **state**

Multiple Steps

Send **reward** and **next state** from a Markovian transition dynamics



Environment



$$r(s, a), s' \sim P(\cdot | s, a)$$

$$s_0 \sim \mu_0, a_0 \sim \pi(s_0), r_0, s_1 \sim P(s_0, a_0), a_1 \sim \pi(s_1), r_1 \dots$$








	Learn from Experience	Generalize	Interactive	Exploration	Credit assignment
Supervised Learning					
Reinforcement Learning					

	Learn from Experience	Generalize	Interactive	Exploration	Credit assignment
Supervised Learning					
Reinforcement Learning					

	Learn from Experience	Generalize	Interactive	Exploration	Credit assignment
Supervised Learning	✓	✓			
Reinforcement Learning	✓	✓			

	Learn from Experience	Generalize	Interactive	Exploration	Credit assignment
Supervised Learning	✓	✓			
Reinforcement Learning	✓	✓	✓		

	Learn from Experience	Generalize	Interactive	Exploration	Credit assignment
Supervised Learning	✓	✓			
Reinforcement Learning	✓	✓	✓	✓	

	Learn from Experience	Generalize	Interactive	Exploration	Credit assignment
Supervised Learning					
Reinforcement Learning					

Infinite horizon Discounted MDP

$$\mathcal{M} = \{S, A, P, r, \mu_0, \gamma\}$$

$$P : S \times A \mapsto \Delta(S), \quad r : S \times A \rightarrow [0,1], \quad \gamma \in [0,1)$$

Infinite horizon Discounted MDP

$$\mathcal{M} = \{S, A, P, r, \mu_0, \gamma\}$$

$$P : S \times A \mapsto \Delta(S), \quad r : S \times A \rightarrow [0,1], \quad \gamma \in [0,1)$$

$$\text{Policy } \pi : S \mapsto \Delta(A)$$

Infinite horizon Discounted MDP

$$\mathcal{M} = \{S, A, P, r, \mu_0, \gamma\}$$

$$P : S \times A \mapsto \Delta(S), \quad r : S \times A \rightarrow [0,1], \quad \gamma \in [0,1)$$

$$\text{Policy } \pi : S \mapsto \Delta(A)$$

$$\text{Value function } V^\pi(s) = \mathbb{E} \left[\sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \mid s_0 = s, a_h \sim \pi(s_h), s_{h+1} \sim P(\cdot \mid s_h, a_h) \right]$$

Infinite horizon Discounted MDP

$$\mathcal{M} = \{S, A, P, r, \mu_0, \gamma\}$$

$$P : S \times A \mapsto \Delta(S), \quad r : S \times A \rightarrow [0,1], \quad \gamma \in [0,1)$$

$$\text{Policy } \pi : S \mapsto \Delta(A)$$

$$\text{Value function } V^\pi(s) = \mathbb{E} \left[\sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \mid s_0 = s, a_h \sim \pi(s_h), s_{h+1} \sim P(\cdot \mid s_h, a_h) \right]$$

$$\text{Q function } Q^\pi(s, a) = \mathbb{E} \left[\sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \mid (s_0, a_0) = (s, a), a_h \sim \pi(s_h), s_{h+1} \sim P(\cdot \mid s_h, a_h) \right]$$

Bellman Equation:

$$V^\pi(s) = \mathbb{E} \left[\sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \mid s_0 = s, a_h \sim \pi(s_h), s_{h+1} \sim P(\cdot \mid s_h, a_h) \right]$$

Bellman Equation:

$$V^\pi(s) = \mathbb{E} \left[\sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \mid s_0 = s, a_h \sim \pi(s_h), s_{h+1} \sim P(\cdot \mid s_h, a_h) \right]$$

$$V^\pi(s) = \mathbb{E}_{a \sim \pi(s)} \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot \mid s, a)} V^\pi(s') \right]$$

Bellman Equation:

$$V^\pi(s) = \mathbb{E} \left[\sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \mid s_0 = s, a_h \sim \pi(s_h), s_{h+1} \sim P(\cdot \mid s_h, a_h) \right]$$

$$V^\pi(s) = \mathbb{E}_{a \sim \pi(s)} \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot \mid s, a)} V^\pi(s') \right]$$

$$Q^\pi(s, a) = \mathbb{E} \left[\sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \mid (s_0, a_0) = (s, a), a_h \sim \pi(s_h), s_{h+1} \sim P(\cdot \mid s_h, a_h) \right]$$

Bellman Equation:

$$V^\pi(s) = \mathbb{E} \left[\sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \mid s_0 = s, a_h \sim \pi(s_h), s_{h+1} \sim P(\cdot \mid s_h, a_h) \right]$$

$$V^\pi(s) = \mathbb{E}_{a \sim \pi(s)} \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot \mid s, a)} V^\pi(s') \right]$$

$$Q^\pi(s, a) = \mathbb{E} \left[\sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \mid (s_0, a_0) = (s, a), a_h \sim \pi(s_h), s_{h+1} \sim P(\cdot \mid s_h, a_h) \right]$$

$$Q^\pi(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot \mid s, a)} V^\pi(s')$$

Outline

 1. Definition of infinite horizon discounted MDPs

2. Bellman Optimality

3. State-action distribution

Optimal Policy

For infinite horizon discounted MDP, there exists a deterministic stationary policy

$$\pi^{\star} : S \mapsto A, \text{ s.t.}, V^{\pi^{\star}}(s) \geq V^{\pi}(s), \forall s, \pi$$

[Puterman 94 chapter 6, also see theorem 1.4 in the RL monograph]

Optimal Policy

For infinite horizon discounted MDP, there exists a deterministic stationary policy

$$\pi^{\star} : S \mapsto A, \text{ s.t.}, V^{\pi^{\star}}(s) \geq V^{\pi}(s), \forall s, \pi$$

[Puterman 94 chapter 6, also see theorem 1.4 in the RL monograph]

We denote $V^{\star} := V^{\pi^{\star}}, Q^{\star} := Q^{\pi^{\star}}$

Optimal Policy

For infinite horizon discounted MDP, there exists a deterministic stationary policy

$$\pi^* : S \mapsto A, \text{ s.t.}, V^{\pi^*}(s) \geq V^\pi(s), \forall s, \pi$$

[Puterman 94 chapter 6, also see theorem 1.4 in the RL monograph]

We denote $V^* := V^{\pi^*}, Q^* := Q^{\pi^*}$

Theorem 1: Bellman Optimality

$$V^*(s) = \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V^*(s') \right]$$

Proof of Bellman Optimality

Theorem 1: Bellman Optimality

$$V^*(s) = \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V^*(s') \right]$$

Proof of Bellman Optimality

Theorem 1: Bellman Optimality

$$V^*(s) = \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V^*(s') \right]$$

Denote $\hat{\pi}(s) := \arg \max_a Q^*(s, a)$, we will prove $V^{\hat{\pi}}(s) = V^*(s), \forall s$

Proof of Bellman Optimality

Theorem 1: Bellman Optimality

$$V^*(s) = \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V^*(s') \right]$$

Denote $\hat{\pi}(s) := \arg \max_a Q^*(s, a)$, we will prove $V^{\hat{\pi}}(s) = V^*(s), \forall s$

$$V^*(s) = r(s, \pi^*(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi^*(s))} V^*(s')$$

Proof of Bellman Optimality

Theorem 1: Bellman Optimality

$$V^*(s) = \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V^*(s') \right]$$

Denote $\hat{\pi}(s) := \arg \max_a Q^*(s, a)$, we will prove $V^{\hat{\pi}}(s) = V^*(s), \forall s$

$$\begin{aligned} V^*(s) &= r(s, \pi^*(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi^*(s))} V^*(s') \\ &\leq \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s') \right] = r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} V^*(s') \end{aligned}$$

Proof of Bellman Optimality

Theorem 1: Bellman Optimality

$$V^*(s) = \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V^*(s') \right]$$

Denote $\hat{\pi}(s) := \arg \max_a Q^*(s, a)$, we will prove $V^{\hat{\pi}}(s) = V^*(s), \forall s$

$$\begin{aligned} V^*(s) &= r(s, \pi^*(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi^*(s))} V^*(s') \\ &\leq \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s') \right] = r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} V^*(s') \\ &= r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} \left[r(s', \pi^*(s')) + \gamma \mathbb{E}_{s'' \sim P(s', \pi^*(s'))} V^*(s'') \right] \end{aligned}$$

Proof of Bellman Optimality

Theorem 1: Bellman Optimality

$$V^*(s) = \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V^*(s') \right]$$

Denote $\hat{\pi}(s) := \arg \max_a Q^*(s, a)$, we will prove $V^{\hat{\pi}}(s) = V^*(s), \forall s$

$$\begin{aligned} V^*(s) &= r(s, \pi^*(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi^*(s))} V^*(s') \\ &\leq \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s') \right] = r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} V^*(s') \\ &= r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} \left[r(s', \pi^*(s')) + \gamma \mathbb{E}_{s'' \sim P(s', \pi^*(s'))} V^*(s'') \right] \\ &\leq r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} \left[r(s', \hat{\pi}(s')) + \gamma \mathbb{E}_{s'' \sim P(s', \hat{\pi}(s'))} V^*(s'') \right] \end{aligned}$$

Proof of Bellman Optimality

Theorem 1: Bellman Optimality

$$V^*(s) = \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V^*(s') \right]$$

Denote $\hat{\pi}(s) := \arg \max_a Q^*(s, a)$, we will prove $V^{\hat{\pi}}(s) = V^*(s), \forall s$

$$\begin{aligned} V^*(s) &= r(s, \pi^*(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi^*(s))} V^*(s') \\ &\leq \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s') \right] = r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} V^*(s') \\ &= r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} \left[r(s', \pi^*(s')) + \gamma \mathbb{E}_{s'' \sim P(s', \pi^*(s'))} V^*(s'') \right] \\ &\leq r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} \left[r(s', \hat{\pi}(s')) + \gamma \mathbb{E}_{s'' \sim P(s', \hat{\pi}(s'))} V^*(s'') \right] \\ &\leq r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} \left[r(s', \hat{\pi}(s')) + \gamma \mathbb{E}_{s'' \sim P(s', \hat{\pi}(s'))} \left[r(s'', \hat{\pi}(s'')) + \gamma \mathbb{E}_{s''' \sim P(s'', \hat{\pi}(s''))} V^*(s''') \right] \right] \end{aligned}$$

Proof of Bellman Optimality

Theorem 1: Bellman Optimality

$$V^*(s) = \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V^*(s') \right]$$

Denote $\hat{\pi}(s) := \arg \max_a Q^*(s, a)$, we will prove $V^{\hat{\pi}}(s) = V^*(s), \forall s$

$$\begin{aligned} V^*(s) &= r(s, \pi^*(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi^*(s))} V^*(s') \\ &\leq \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s') \right] = r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} V^*(s') \\ &= r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} \left[r(s', \pi^*(s')) + \gamma \mathbb{E}_{s'' \sim P(s', \pi^*(s'))} V^*(s'') \right] \\ &\leq r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} \left[r(s', \hat{\pi}(s')) + \gamma \mathbb{E}_{s'' \sim P(s', \hat{\pi}(s'))} V^*(s'') \right] \\ &\leq r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} \left[r(s', \hat{\pi}(s')) + \gamma \mathbb{E}_{s'' \sim P(s', \hat{\pi}(s'))} \left[r(s'', \hat{\pi}(s'')) + \gamma \mathbb{E}_{s''' \sim P(s'', \hat{\pi}(s''))} V^*(s''') \right] \right] \\ &\leq \mathbb{E} \left[r(s, \hat{\pi}(s)) + \gamma r(s', \hat{\pi}(s')) + \dots \right] = V^{\hat{\pi}}(s) \end{aligned}$$

Proof of Bellman Optimality

Theorem 1: Bellman Optimality

$$V^*(s) = \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V^*(s') \right]$$

Denote $\hat{\pi}(s) := \arg \max_a Q^*(s, a)$, we just proved $V^{\hat{\pi}}(s) = V^*(s), \forall s$

Proof of Bellman Optimality

Theorem 1: Bellman Optimality

$$V^*(s) = \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V^*(s') \right]$$

Denote $\hat{\pi}(s) := \arg \max_a Q^*(s, a)$, we just proved $V^{\hat{\pi}}(s) = V^*(s), \forall s$

This implies that $\arg \max_a Q^*(s, a)$ is an optimal policy

Proof of Bellman Optimality

Theorem 2:

For any $V : S \rightarrow \mathbb{R}$, if $V(s) = \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V(s') \right]$ for all s ,
then $V(s) = V^*(s), \forall s$

Proof of Bellman Optimality

Theorem 2:

For any $V : S \rightarrow \mathbb{R}$, if $V(s) = \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V(s') \right]$ for all s ,
then $V(s) = V^*(s), \forall s$

$$|V(s) - V^*(s)| = \left| \max_a (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V(s')) - \max_a (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s')) \right|$$

Proof of Bellman Optimality

Theorem 2:

For any $V : S \rightarrow \mathbb{R}$, if $V(s) = \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V(s') \right]$ for all s ,
then $V(s) = V^*(s), \forall s$

$$\begin{aligned} |V(s) - V^*(s)| &= \left| \max_a (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V(s')) - \max_a (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s')) \right| \\ &\leq \max_a \left| (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V(s')) - (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s')) \right| \end{aligned}$$

Proof of Bellman Optimality

Theorem 2:

For any $V : S \rightarrow \mathbb{R}$, if $V(s) = \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V(s') \right]$ for all s ,
then $V(s) = V^*(s), \forall s$

$$\begin{aligned} |V(s) - V^*(s)| &= \left| \max_a (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V(s')) - \max_a (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s')) \right| \\ &\leq \max_a \left| (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V(s')) - (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s')) \right| \\ &\leq \max_a \gamma \mathbb{E}_{s' \sim P(s, a)} |V(s') - V^*(s')| \end{aligned}$$

Proof of Bellman Optimality

Theorem 2:

For any $V : S \rightarrow \mathbb{R}$, if $V(s) = \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V(s') \right]$ for all s ,
then $V(s) = V^*(s), \forall s$

$$\begin{aligned} |V(s) - V^*(s)| &= \left| \max_a (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V(s')) - \max_a (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s')) \right| \\ &\leq \max_a \left| (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V(s')) - (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s')) \right| \\ &\leq \max_a \gamma \mathbb{E}_{s' \sim P(s, a)} \left| V(s') - V^*(s') \right| \\ &\leq \max_a \gamma \mathbb{E}_{s' \sim P(s, a)} \left(\max_{a'} \gamma \mathbb{E}_{s'' \sim P(s', a')} \left| V(s'') - V^*(s'') \right| \right) \end{aligned}$$

Proof of Bellman Optimality

Theorem 2:

For any $V : S \rightarrow \mathbb{R}$, if $V(s) = \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V(s') \right]$ for all s ,
then $V(s) = V^*(s), \forall s$

$$\begin{aligned} |V(s) - V^*(s)| &= \left| \max_a (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V(s')) - \max_a (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s')) \right| \\ &\leq \max_a \left| (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V(s')) - (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s')) \right| \\ &\leq \max_a \gamma \mathbb{E}_{s' \sim P(s, a)} \left| V(s') - V^*(s') \right| \\ &\leq \max_a \gamma \mathbb{E}_{s' \sim P(s, a)} \left(\max_{a'} \gamma \mathbb{E}_{s'' \sim P(s', a')} \left| V(s'') - V^*(s'') \right| \right) \\ &\leq \max_{a_1, a_2, \dots, a_{k-1}} \gamma^k \mathbb{E}_{s_k} |V(s_k) - V^*(s_k)| \end{aligned}$$

Outline

 1. Definition of infinite horizon discounted MDPs

 2. Bellman Optimality

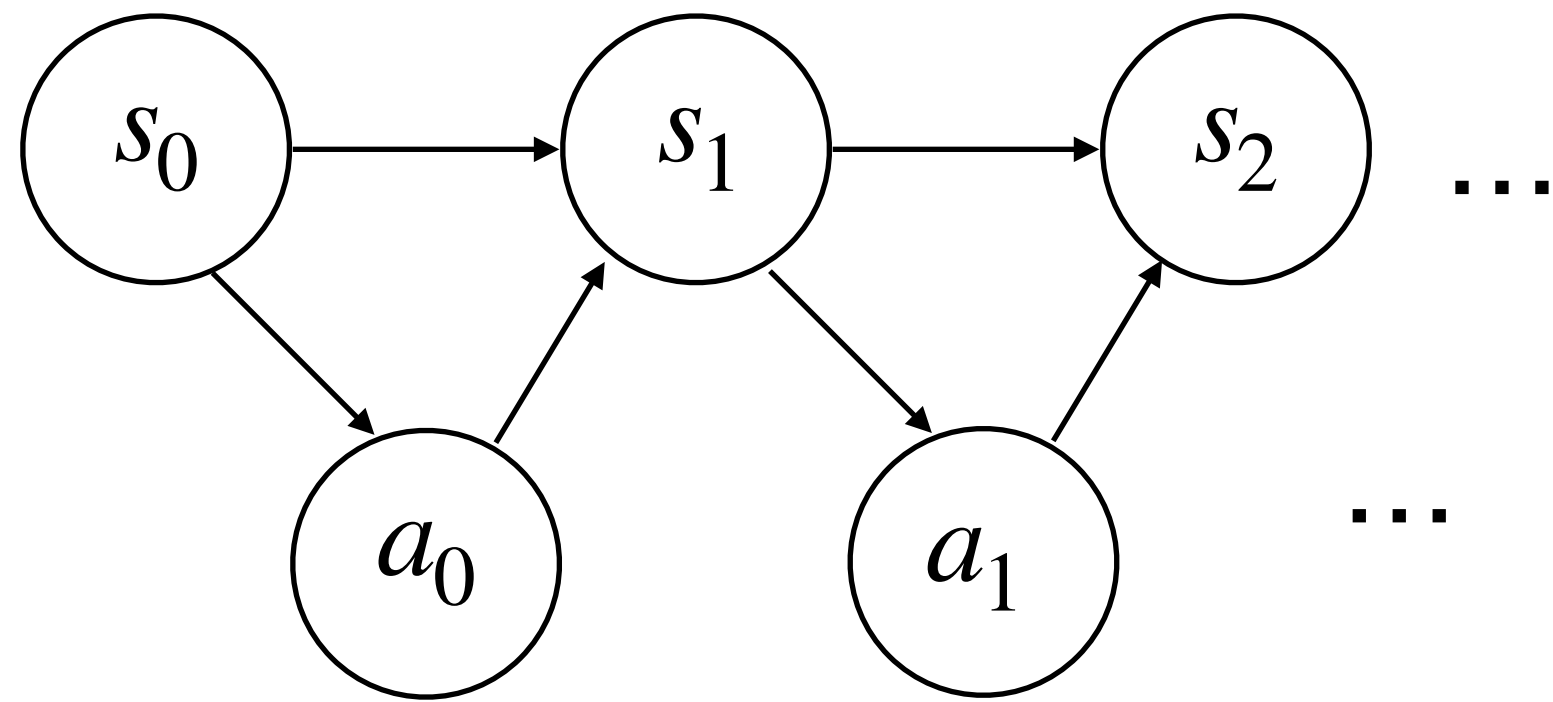
3. State-action distribution

Trajectory distribution and state-action distribution

Q: what is the probability of π generating trajectory $\tau = \{s_0, a_0, s_1, a_1, \dots, s_h, a_h\}$?

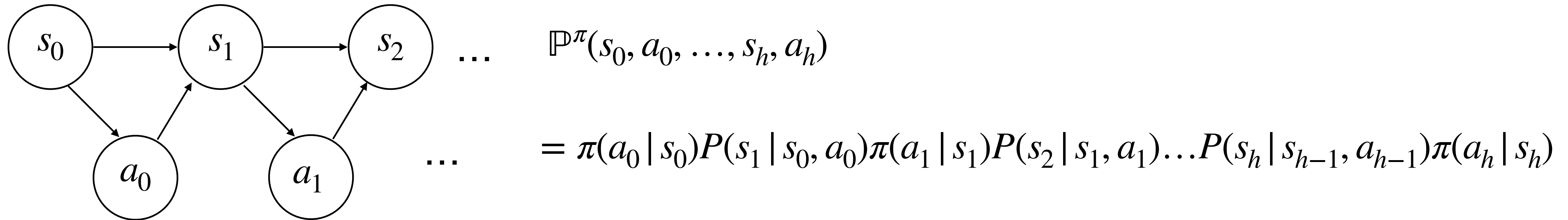
Trajectory distribution and state-action distribution

Q: what is the probability of π generating trajectory $\tau = \{s_0, a_0, s_1, a_1, \dots, s_h, a_h\}$?



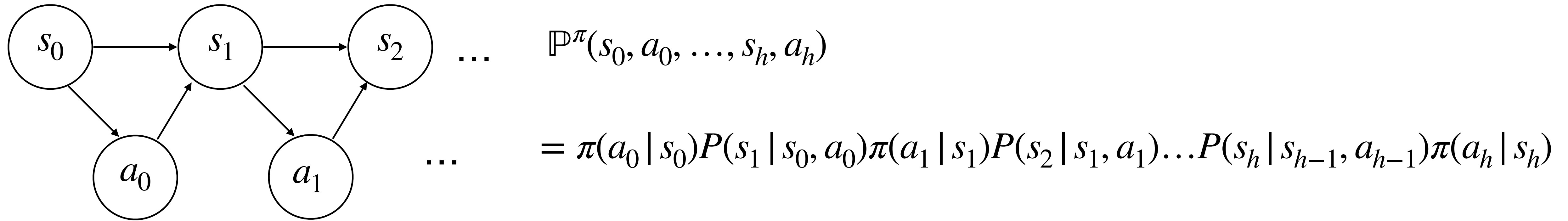
Trajectory distribution and state-action distribution

Q: what is the probability of π generating trajectory $\tau = \{s_0, a_0, s_1, a_1, \dots, s_h, a_h\}$?



Trajectory distribution and state-action distribution

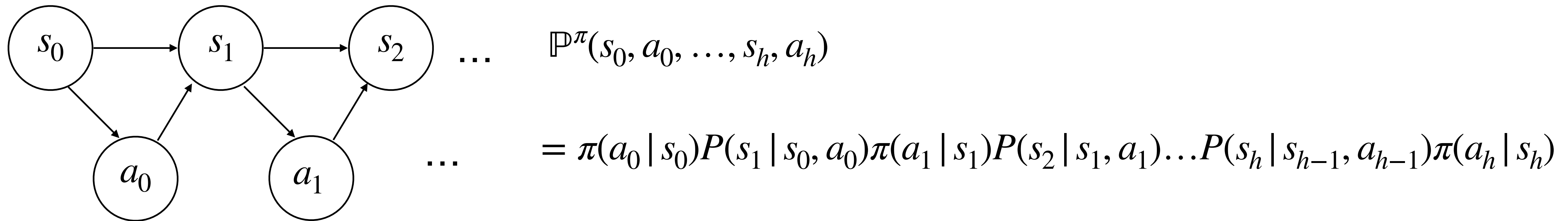
Q: what is the probability of π generating trajectory $\tau = \{s_0, a_0, s_1, a_1, \dots, s_h, a_h\}$?



Q: what's the probability of π visiting state (s,a) at time step h ?

Trajectory distribution and state-action distribution

Q: what is the probability of π generating trajectory $\tau = \{s_0, a_0, s_1, a_1, \dots, s_h, a_h\}$?



Q: what's the probability of π visiting state (s, a) at time step h ?

$$\mathbb{P}_h^\pi(s, a; s_0) = \sum_{a_0, s_1, a_1, \dots, s_{h-1}, a_{h-1}} \mathbb{P}^\pi(s_0, a_0, \dots, s_{h-1}, a_{h-1} | s_h = s, a_h = a)$$

State action occupancy measure

$\mathbb{P}_h(s, a; s_0, \pi)$: probability of π visiting (s, a) at time step $h \in \mathbb{N}$, starting at s_0

$$d_{s_0}^\pi(s, a) = (1 - \gamma) \sum_{h=0}^{\infty} \gamma^h \mathbb{P}_h(s, a; s_0, \pi)$$

$$V^\pi(s_0) = \frac{1}{1 - \gamma} \sum_{s,a} d_{s_0}^\pi(s, a) r(s, a)$$

Summary for today

Key definitions: MDPs, Value / Q functions, State-action distribution

Key property: Bellman optimality (the two theorems and their proofs)