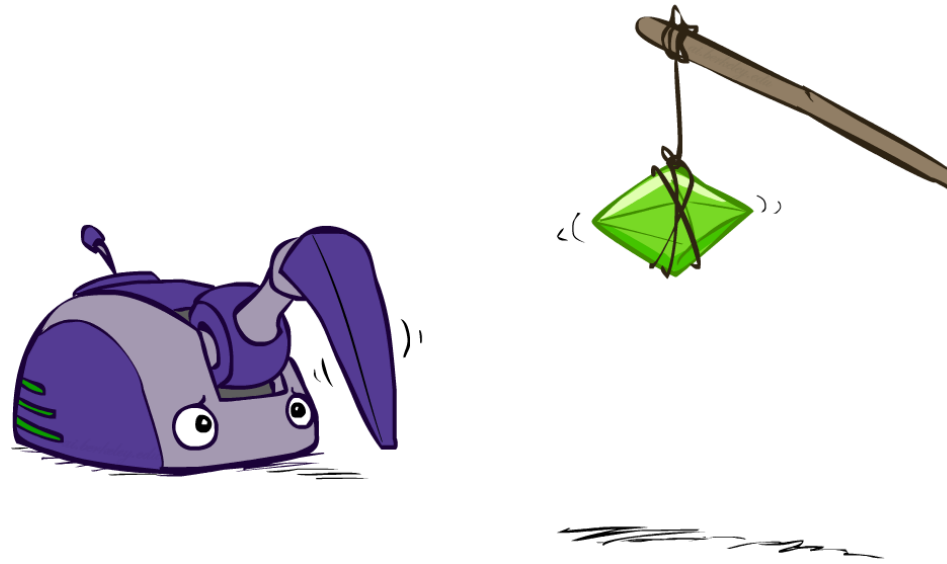


# 15-281 AI: Representation and Problem Solving

## Reinforcement Learning I



Instructors: Nihar Shah and Tuomas Sandholm

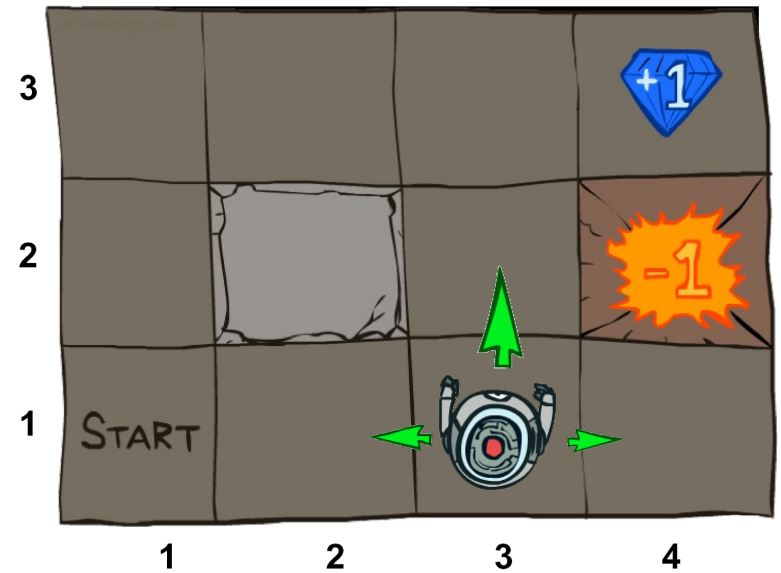
Carnegie Mellon University

Slide credits: CMU AI and ai.berkeley.edu

# Recall Markov Decision Processes

- An MDP is defined by:
  - A **set of states**  $s \in S$
  - A **set of actions**  $a \in A$
  - A **transition function**  $T(s, a, s')$ 
    - Probability that  $a$  from  $s$  leads to  $s'$ , i.e.,  $P(s' | s, a)$
  - A **reward function**  $R(s, a, s')$

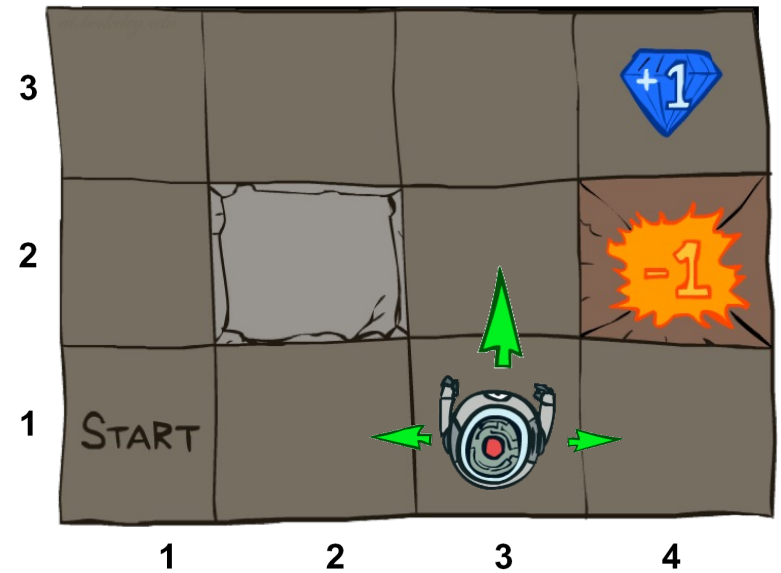
All of this information is known beforehand



# Reinforcement Learning (RL)

- RL is defined by:
  - A **set of states**  $s \in S$
  - A **set of actions**  $a \in A$
  - A **transition function**  $T(s, a, s')$ 
    - Probability that  $a$  from  $s$  leads to  $s'$ , i.e.,  $P(s' | s, a)$
  - A **reward function**  $R(s, a, s')$

Then what is the difference with MDPs?



# MDPs vs. Reinforcement Learning

---

## MDPs

- You know the rewards and transitions (i.e., the model of the world) beforehand
- E.g., a prior rover etc. mapped the entire terrain which you can use now



## Reinforcement Learning

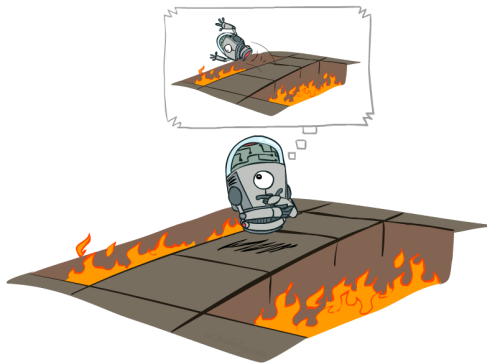
- You don't know rewards and transitions beforehand
- E.g., the robot is deployed in an environment that is completely unseen before

# MDPs vs. Reinforcement Learning

---

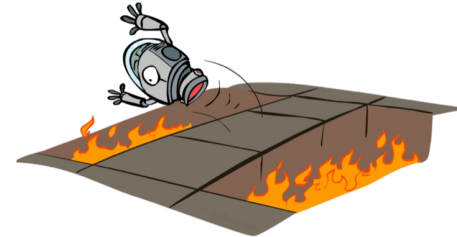
## MDPs

- You know the rewards and transitions (i.e., the model of the world) beforehand
- Find policy to maximize total reward
- **Run our MDP solvers offline**
- **Only deploy the optimal policy**



## Reinforcement Learning

- You don't know rewards and transitions beforehand
- Still assume the world is an MDP
- Do not know how the world works
- Find policy to maximize total reward
- **Need to take actions in the world to understand the world**



# DeepMind Atari (©Two Minute Lectures)

---



# Example: Learning to Walk

---



Initial

[Kohl and Stone, ICRA 2004]

[Video: AIBO WALK – initial]

# Example: Learning to Walk

---



[Kohl and Stone, ICRA 2004]

Training

[Video: AIBO WALK – traini



# Example: Learning to Walk

---



Finished

[Kohl and Stone, ICRA 2004]

[Video: AIBO WALK – finish]

# Reinforcement Learning

---

We don't know  $T$  and  $R$ ! How do we find policies?

We need to estimate quantities by trying out

At time step  $t$ ,

- Take action  $a_t$
- End up in new state  $s_{t+1}$
- Observe reward =  $R(s_t, a_t, s_{t+1})$

Data:  $s_0, a_1, R_1, s_1, a_2, R_2, s_2, a_3, R_3, s_3, \dots$

# Passive vs Active RL

---

- Passive RL
  - Suppose policy  $\pi$  being used is already given to you
  - Only worry about how to learn from experience
  - Agent does **not** control the policy  $\pi$  of taking actions
- Active RL
  - Agent has to decide how to collect experience

# Today: Passive RL

---

- Let  $\pi$  denote the policy being used
- One may like to learn the value obtained from this policy:
  - $V^\pi(s)$  = expected reward starting at state  $s$  and always following  $\pi$
  - $Q^\pi(s, a)$  = expected reward starting at  $s$ , taking action  $a$  now, and following  $\pi$  thereafter
  - One may like to learn the optimal policy:  $\pi^*$  or  $V^*(s)$  or  $Q^*(s, a)$

Data:  $s_0, a_1, R_1, s_1, a_2, R_2, s_2, a_3, R_3, s_3, \dots$

- Data collected via a policy  $\pi$

# Today: Policy evaluation

---

- For the policy  $\pi$  being used in the passive RL setting, how to estimate  $V^\pi$ ?
- (Next lecture: How do we compute an optimal policy)

# Monte-Carlo estimation

---

*Collect a bunch of samples and take their average value*

- Suppose we want to compute the expected age of people in the US
- **Strategy:** Sample  $N$  people at random. Suppose their ages are  $a_1, a_2, \dots, a_N$ . Now “estimate” the expected age to be

$$E[A] \approx \frac{1}{N} (a_1 + a_2 + \dots + a_N)$$

# Reinforcement learning approaches

---

- **Model based RL**

- First estimate  $T(s, a, s')$  and  $R(s, a, s')$
- Then use methods from MDP



- Data:  $s_0, a_1, R_1, s_1, a_2, R_2, s_2, a_3, R_3, s_3, \dots$

- $\hat{T}(s, a, s') = \frac{\# \text{ times } (s, a, s') \text{ occurs}}{\# \text{ times } (s, a) \text{ occurs}}$

- $\hat{R}(s, a, s') = \text{mean}(\text{rewards obtained from state } s \text{ taking action } a \text{ and moving to } s')$

# Reinforcement learning approaches

---

- **Model based RL**

- First estimate  $T(s, a, s')$  and  $R(s, a, s')$
- Then use methods from MDP

- **Model free RL**

- Could be wasteful to estimate  $T(s, a, s')$  and  $R(s, a, s')$
- Estimate  $V^\pi$  and  $Q^\pi$  directly
- Estimate  $Q^*$  directly (policy extraction is trivial)



# Policy evaluation

---

- For the policy  $\pi$  being used in the passive RL setting, how to estimate  $V^\pi$ ?

Rest of lecture: Three methods (progressively improving):

- Direct policy evaluation
- “Better” policy evaluation
- Temporal difference learning

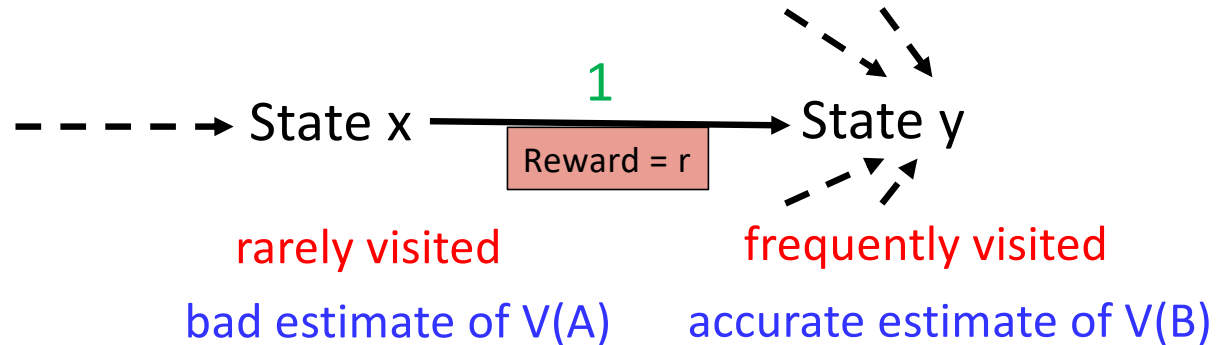
# Direct policy evaluation

---

- Estimate  $V^\pi$  as follows
- Every time you visit a state, write down what the total reward across time turned out to be
- Average those samples
- *Monte-carlo estimates using samples of utility*

# Direct policy evaluation

- **Problem:** This does **not** utilize the MDP structure, and can be wasteful



But, if we were to exploit the MDP structure:  $V^\pi(x) = r + \gamma V^\pi(y)$

# “Better” policy evaluation

- Recall value iteration (for given policy  $\pi$ )

$$V_{k+1}^{\pi}(s) \leftarrow \sum_{s'} T(s, \pi(s), s') [R(s, \pi(s), s') + \gamma V_k^{\pi}(s')]$$

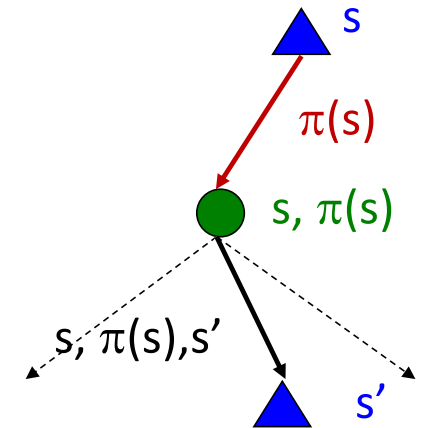
- Can we do this **without T, R** but just **via samples**?
  - Consider each time you are in state  $s$  and take action  $\pi(s)$

$$sample_1 = R(s, \pi(s), s'_1) + \gamma V_k^{\pi}(s'_1)$$

$$sample_2 = R(s, \pi(s), s'_2) + \gamma V_k^{\pi}(s'_2)$$

$$\dots$$
$$sample_n = R(s, \pi(s), s'_n) + \gamma V_k^{\pi}(s'_n)$$

$$V_{k+1}^{\pi}(s) \leftarrow \frac{1}{n} \sum_i sample_i$$



# Challenges with “better” policy estimation

---

$$\text{sample}_1 = R(s, \pi(s), s'_1) + \gamma V_k^\pi(s'_1)$$

$$\dots$$
$$\text{sample}_n = R(s, \pi(s), s'_n) + \gamma V_k^\pi(s'_n)$$

$$V_{k+1}^\pi(s) \leftarrow \frac{1}{n} \sum_i \text{sample}_i$$

- Need to wait for  $n$  samples of  $(s, \pi(s))$  for the update
- Any update to a different state  $V^\pi(s')$  in the meantime will use a stale estimate of  $V^\pi(s)$
- Update for  $V^\pi(s)$  may be using stale estimates of  $V^\pi(s')$

# Reinterpreting sample average

---

$$V_{k+1}^{\pi}(s) \leftarrow \frac{1}{n} \sum_i \text{sample}_i$$

- Suppose we have average over  $n$  samples, and see a new sample, how do we update?

$$V_{k+1}^{\pi}(s) = \frac{n}{n+1} V_{k+1}^{\pi}(s) + \frac{1}{n+1} \text{sample}$$

- Interpolating between current estimate and new sample

# Temporal Difference learning

---

- **Main idea:** Update  $V(s)$  each time we experience a transition  $(s, a, s')$

**Sample of  $V^\pi(s)$ :**  $sample = R + \gamma V^\pi(s')$

**Update  $V^\pi(s)$ :**  $V^\pi(s) \leftarrow (1 - \alpha)V^\pi(s) + (\alpha)sample$

- Decreasing learning rate ( $\alpha$ ) towards zero leads to convergence