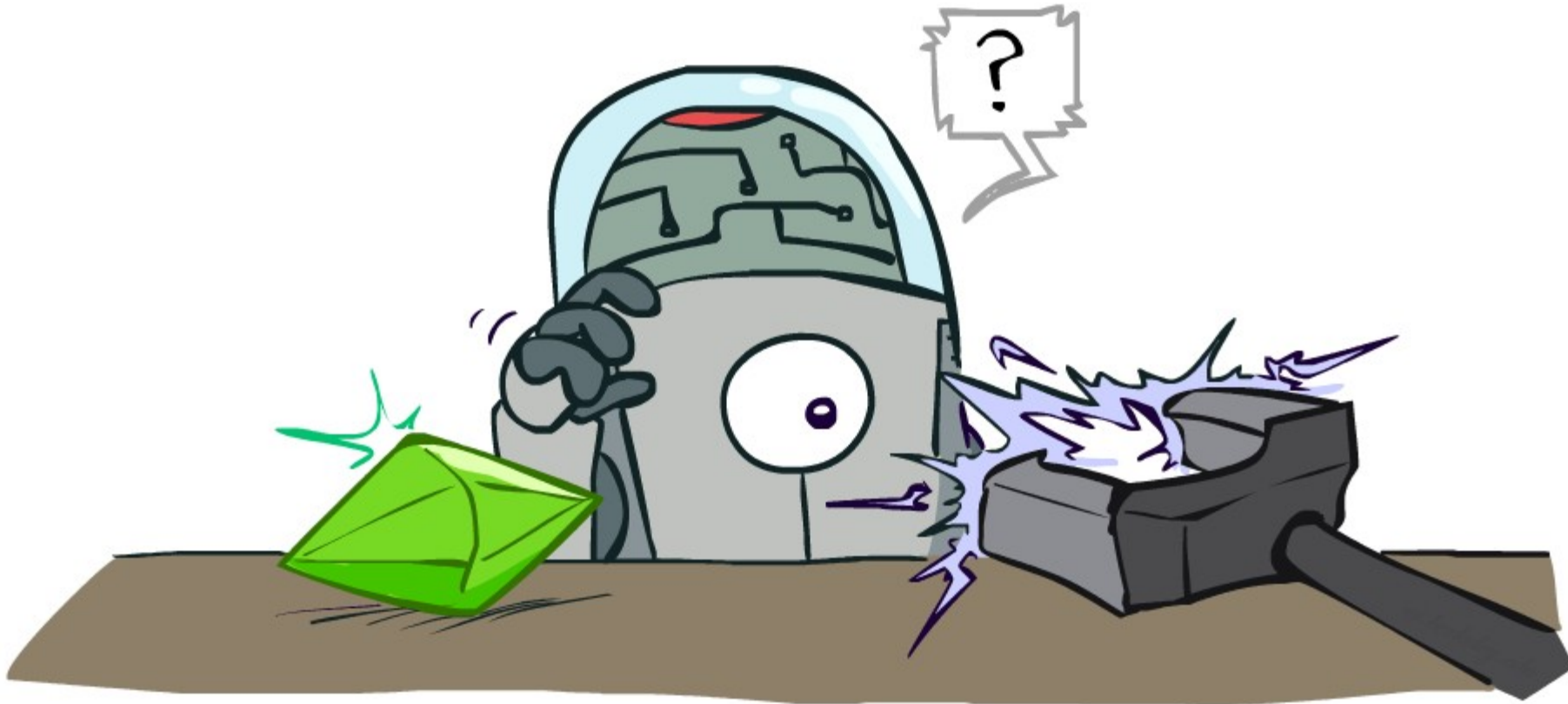


CS 343: Artificial Intelligence

Reinforcement Learning



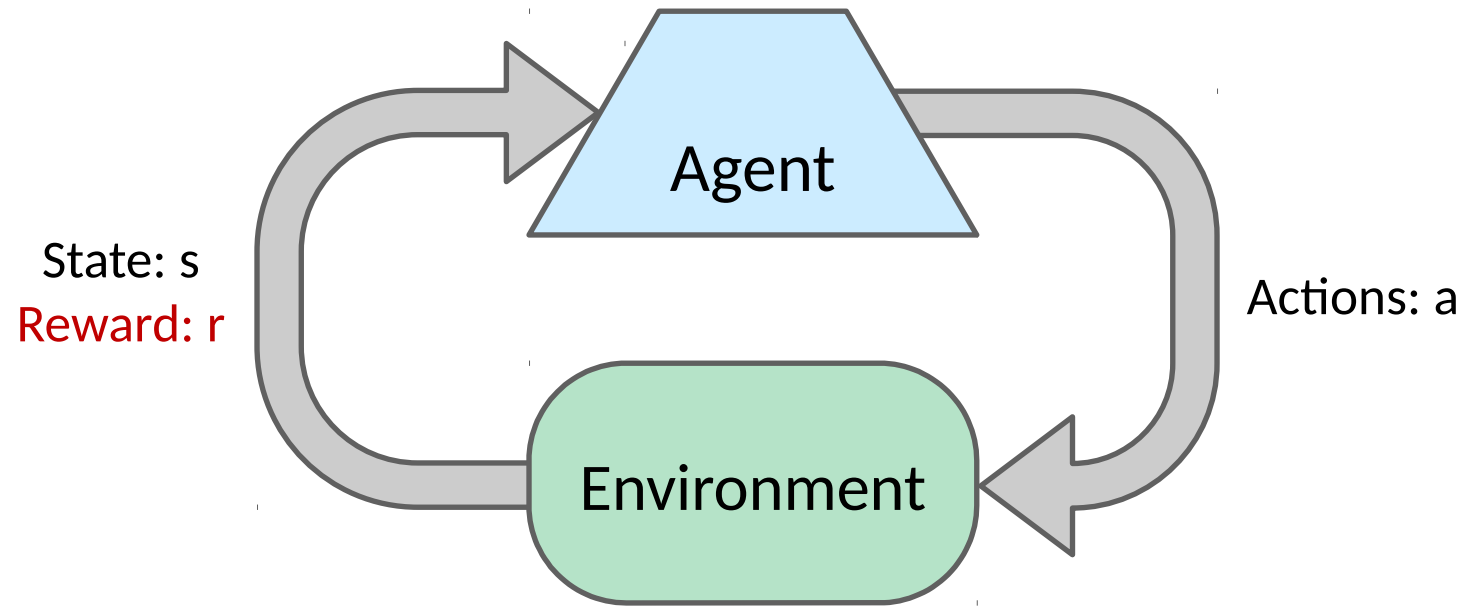
Profs. Peter Stone and Yuke Zhu

The University of Texas at Austin

Good morning colleagues!

- Past due:
 - HW1-3: Search, CSPs, Games
 - 7 reading responses: AI100 report; 6 Textbook readings
 - P0,1: tutorial, Search
- Upcoming EdX Homeworks
 - HW4: MDPs - due Monday 3/8 at 11:59 pm
 - HW5: RL – due Monday 3/22 at 11:59 pm
 - **HW6: Bayes Nets – due Monday 4/5 at 11:59 pm**
- Upcoming programming projects
 - P2: Games – due Wednesday 3/3 at 11:59pm
 - **P3: RL – due Wednesday 3/31 at 11:59pm**
- Readings: Bayes Nets – Due Monday 3/8 at 9:30am
- **Midterm – end of week after spring break (3/25 or so)**
 - Material up through and including Bayes Nets

Reinforcement Learning



- Basic idea:

- Receive feedback in the form of **rewards**
- Agent's utility is defined by the reward function
- Must (learn to) act so as to **maximize expected rewards**
- All learning is based on observed samples of outcomes!

Example: Learning to Walk



Initial



A Learning Trial



After Learning [1K Trials]

Example: Learning to Walk



Initial

Example: Learning to Walk



Training

Example: Learning to Walk



Finished

Some of Your Questions

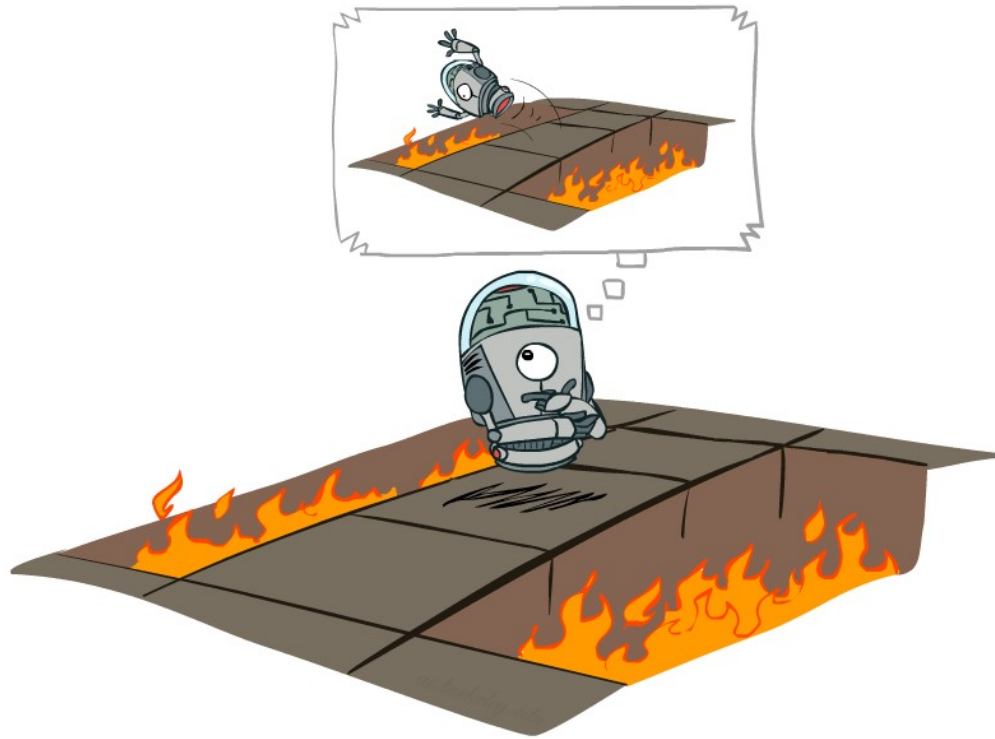
- Difference between MDPs and RL
- Q-learning vs. SARSA
 - What makes Q-learning “off-policy”?
- Model-based vs. Model-free – which is better?
- Passive learning vs. active learning – when would you use each?
 - TD vs. Q-learning
- When should the agent stop learning?
- How does RL relate to “machine learning” and “deep learning”? (Vishal Tak)
- How do you know when your model is accurate enough to start using it? (Aditya Gupta)
- After learning in one environment, does an RL agent work in another? (Rudraksh Garg)
- Are there methods between Monte Carlo and TD that update after arbitrary steps? (Conrad Li)
- Is RL possible without rewards? (Ramya Prasad)
- Do the algorithms still work if there’s more than one agent? (Michael Rodriguez-Labarca)
- How much has RL progressed since the book was written 10 years ago? (Ethan Houston)
 - Was more powerful computing necessary for these advances? (Dale Kang)

Reinforcement Learning

- Still assume a Markov decision process (MDP):
 - A set of states $s \in S$
 - A set of actions (per state) A
 - A model $T(s,a,s')$
 - A reward function $R(s,a,s')$
- Still looking for a policy $\pi(s)$
- New twist: don't know T or R
 - I.e. we don't know which states are good or what the actions do
 -



Offline (MDPs) vs. Online (RL)



Offline Solution



Online Learning

Test Your Understanding

- MDPs and RL
- Practice problem in breakout rooms
- Work for a couple of minutes independently, but then quickly start comparing progress – even if you're not done yet.

Model-Based Learning

- **Model-Based Idea:**
 - Learn an approximate model based on experiences
 - Solve for values as if the learned model were correct
- **Step 1: Learn empirical MDP model**
 - Count outcomes s' for each s, a
 - Normalize to give an estimate of $\hat{T}(s, a, s')$
 - Discover each $\hat{R}(s, a, s')$ when we experience (s, a, s')
- **Step 2: Solve the learned MDP**
 -



Example: Expected Age

Goal: Compute expected age of CS 343 students

Known $P(A)$

$$E[A] = \sum_a P(a) \cdot a = 0.35 \times 20 + \dots$$

Without $P(A)$, instead collect samples $[a_1, a_2, \dots, a_N]$

Unknown $P(A)$: “Model Based”

$$\hat{P}(a) = \frac{\text{num}(a)}{N}$$
$$E[A] \approx \sum_a \hat{P}(a) \cdot a$$

Why does this work? Because eventually you learn the right model.

Unknown $P(A)$: “Model Free”

$$E[A] \approx \frac{1}{N} \sum_i a_i$$

Why does this work? Because samples appear with the right frequencies.

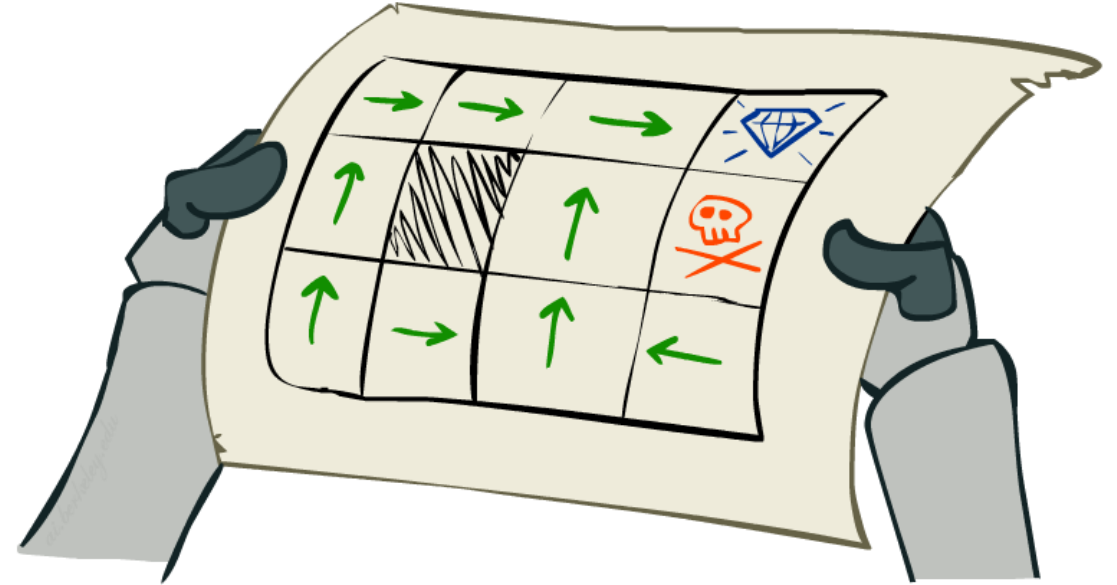
Passive Reinforcement Learning

- Simplified task: policy evaluation

- Input: a fixed policy $\pi(s)$
- You don't know the transitions $T(s,a,s')$
- You don't know the rewards $R(s,a,s')$
- Goal: learn the state values

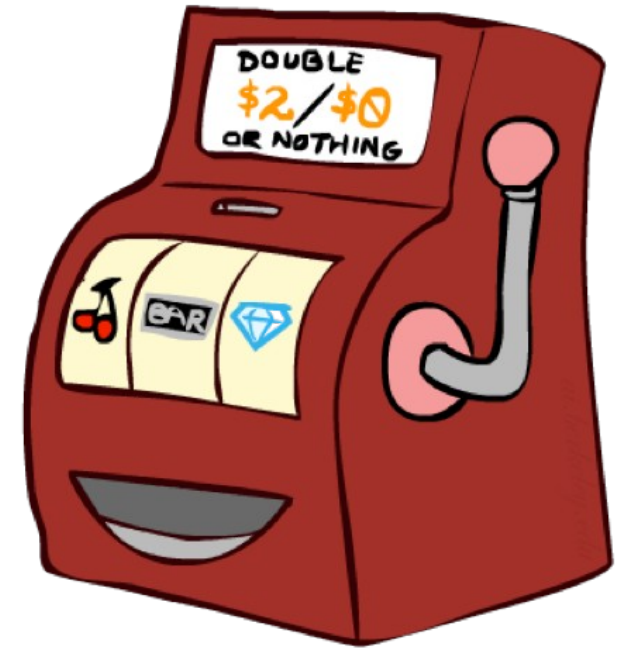
- In this case:

- Learner is “along for the ride”
- No choice about what actions to take
- Just execute the policy and learn from experience
- This is NOT offline planning! You actually take actions in the world.



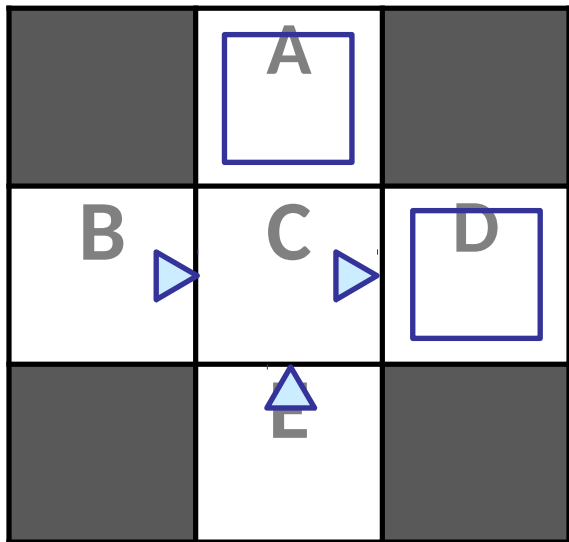
Direct Evaluation

- Goal: Compute values for each state under π
- Idea: Average together observed sample values
 - Act according to π
 - Every time you visit a state, write down what the sum of discounted rewards turned out to be
 - Average those samples
- This is called direct evaluation



Example: Direct Evaluation

Input Policy π



Assume: $\gamma = 1$

Observed Episodes (Training)

Episode 1

B, east, C, -1
C, east, D, -1
D, exit, x, +10

Episode 2

B, east, C, -1
C, east, D, -1
D, exit, x, +10

Episode 3

E, north, C, -1
C, east, D, -1
D, exit, x, +10

Episode 4

E, north, C, -1
C, east, A, -1
A, exit, x, -10

Output Values

	A -10	
B +8	C +4	D +10
	E -2	

Problems with Direct Evaluation

- What's good about direct evaluation?
 - It's easy to understand
 - It doesn't require any knowledge of T , R
 - It eventually computes the correct average values, using just sample transitions
- What bad about it?
 - It wastes information about state connections
 - Each state must be learned separately
 -

Output Values

	A-10	
B+8	C+4	D+10
	E-2	

If B and E both go to C under this policy, how can their values be different?

Why Not Use Policy Evaluation?

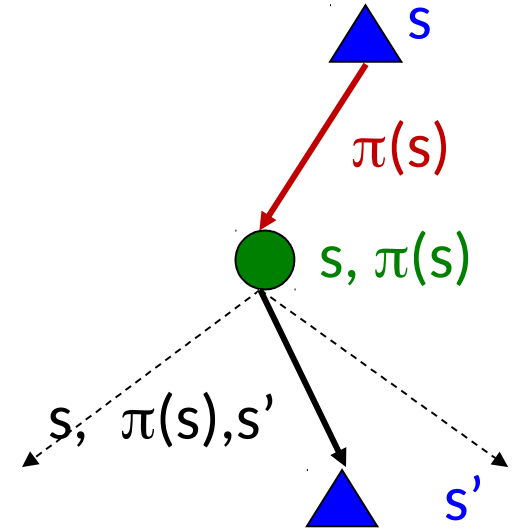
- Simplified Bellman updates calculate V for a fixed policy:

- Each round, replace V with a one-step-look-ahead layer over V

$$V_0^\pi(s) = 0$$

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

- This approach fully exploited the connections between the states
- Unfortunately, we need T and R to do it!



- Key question: how can we do this update to V without knowing T and R ?

-

Sample-Based Policy Evaluation?

- We want to improve our estimate of V by computing these averages:

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

- Idea: Take samples of outcomes s' (by doing the action!) and average

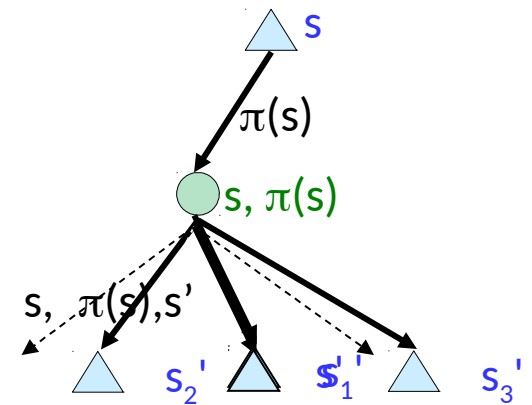
$$sample_1 = R(s, \pi(s), s'_1) + \gamma V_k^{\pi}(s'_1)$$

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

...

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

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

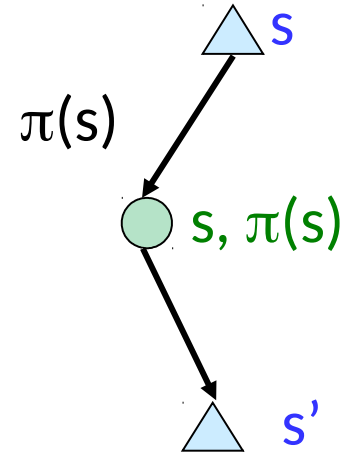


Almost! But we can't
rewind time to get sample
after sample from state s .

Temporal Difference Learning

- Big idea: learn from every experience!

- Update $V(s)$ each time we experience a transition (s, a, s', r)
- Likely outcomes s' will contribute updates more often



- Temporal difference learning of values

- ✦ Policy still fixed, still doing evaluation!
- ✦

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

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

Same update: $V^\pi(s) \leftarrow V^\pi(s) + \alpha(sample - V^\pi(s))$

Exponential Moving Average

- Exponential moving average

- The running interpolation update: $\bar{x}_n = (1 - \alpha) \cdot \bar{x}_{n-1} + \alpha \cdot x_n$

- Makes recent samples more important:

$$\bar{x}_n = \frac{x_n + (1 - \alpha) \cdot x_{n-1} + (1 - \alpha)^2 \cdot x_{n-2} + \dots}{1 + (1 - \alpha) + (1 - \alpha)^2 + \dots}$$

- Forgets about the past (distant past values were wrong anyway)

-