

CS 188: Artificial Intelligence

Reinforcement Learning



Instructor: Marco Alvarez
University of Rhode Island

[These slides were created by Dan Klein and Pieter Abbeel for CS188 Intro to AI at UC Berkeley. All CS188 materials are available at <http://ai.berkeley.edu>.]

1

Double Bandits

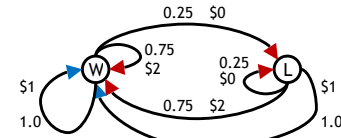


2

Double-Bandit MDP

- Actions: Blue, Red
- States: Win, Lose

No discount
100 time steps
Both states have the same value

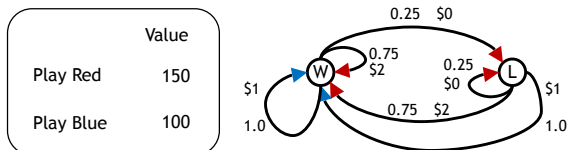


3

Offline Planning

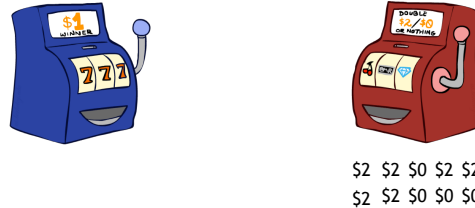
- Solving MDPs is offline planning
 - You determine all quantities through computation
 - You need to know the details of the MDP
 - You do not actually play the game!

No discount
100 time steps
Both states have the same value



4

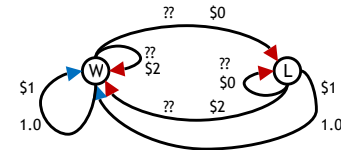
Let's Play!



5

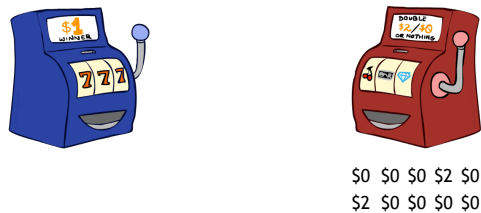
Online Planning

- Rules changed! Red's win chance is different.



6

Let's Play!



7

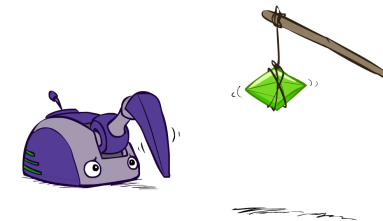
What Just Happened?

- That wasn't planning, it was learning!
 - Specifically, reinforcement learning
 - There was an MDP, but you couldn't solve it with just computation
 - You needed to actually act to figure it out
- Important ideas in reinforcement learning that came up
 - Exploration: you have to try unknown actions to get information
 - Exploitation: eventually, you have to use what you know
 - Regret: even if you learn intelligently, you make mistakes
 - Sampling: because of chance, you have to try things repeatedly
 - Difficulty: learning can be much harder than solving a known MDP



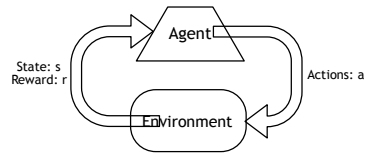
8

Reinforcement Learning



9

Reinforcement Learning

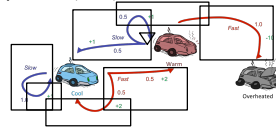


- Basic idea:
 - Receive feedback in the form of rewards
 - Agent's utility is defined by the reward function
 - Must (learn to) act so as to maximize expected rewards
 - All learning is based on observed samples of outcomes!

10

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
 - Must actually try actions and states out to learn



11

Example: Learning to Walk



[Kohl and Stone, ICRA 2004]

Initial

[Video: AIBO WALK - initial]

12

Example: Learning to Walk



[Kohl and Stone, ICRA 2004]

Training

[Video: AIBO WALK - training]

13

Example: Learning to Walk



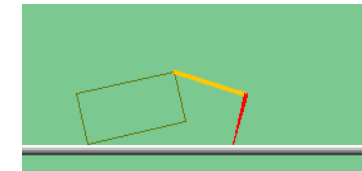
[Kohl and Stone, ICRA 2004]

Finished

[Video: AIBO WALK - finished]

14

The Crawler!



[Demo: Crawler Bot (L10D1)]

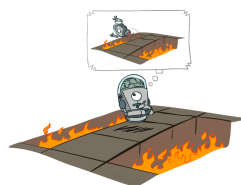
15

Video of Demo Crawler Bot



16

Offline (MDPs) vs. Online (RL)



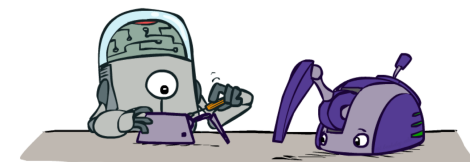
Offline Solution



Online Learning

17

Model-Based Learning



18

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
 - For example, use value iteration, as before



19

Example: Model-Based Learning

Input Policy π	Observed Episodes (Training)		Learned Model									
<table> <tr><td></td><td>A</td><td></td></tr> <tr><td>B</td><td>C</td><td>D</td></tr> <tr><td></td><td>E</td><td></td></tr> </table>		A		B	C	D		E		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	$\hat{T}(s, a, s')$ $\hat{T}(B, \text{east}, C) = 1.00$ $\hat{T}(C, \text{east}, D) = 0.75$ $\hat{T}(C, \text{east}, A) = 0.25$...
	A											
B	C	D										
	E											
	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	$\hat{R}(s, a, s')$ $\hat{R}(B, \text{east}, C) = -1$ $\hat{R}(C, \text{east}, D) = -1$ $\hat{R}(D, \text{exit}, x) = +10$...									
Assume: $\gamma = 1$												

Assume: $\gamma = 1$

20

Example: Expected Age

Goal: Compute expected age of students

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

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

Unknown $P(A)$: "Model Based"

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

$$\hat{P}(a) = \frac{\text{num}(a)}{N}$$

$$E[A] \approx \sum_a \hat{P}(a) \cdot a$$

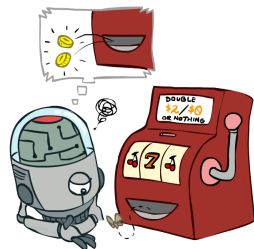
Unknown $P(A)$: "Model Free"

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

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

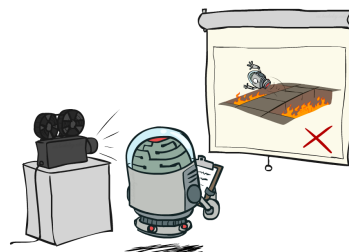
21

Model-Free Learning



22

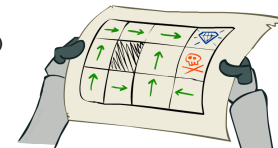
Passive Reinforcement Learning



23

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.



24

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



25

Example: Direct Evaluation

Input Policy π	Observed Episodes (Training)		Output Values																		
<table><tr><td></td><td>A</td><td></td></tr><tr><td>B</td><td>C</td><td>D</td></tr><tr><td></td><td>E</td><td></td></tr></table>		A		B	C	D		E		<div>Episode 1</div> <div>B, east, C, -1 C, east, D, -1 D, exit, x, +10</div>	<div>Episode 2</div> <div>B, east, C, -1 C, east, D, -1 D, exit, x, +10</div>	<table><tr><td></td><td>-10</td><td></td></tr><tr><td>B +8</td><td>C +4</td><td>D +1</td></tr><tr><td></td><td>E -2</td><td></td></tr></table>		-10		B +8	C +4	D +1		E -2	
	A																				
B	C	D																			
	E																				
	-10																				
B +8	C +4	D +1																			
	E -2																				
	<div>Episode 3</div> <div>E, north, C, -1 C, east, D, -1 D, exit, x, +10</div>	<div>Episode 4</div> <div>E, north, C, -1 C, east, A, -1 A, exit, x, -10</div>																			
Assume: $\gamma = 1$																					

Assume: $\gamma = 1$

26

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
 - So, it takes a long time to learn

Output Values

	-10	
+8	C	+10
	E	-2

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

27

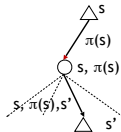
Why Not Use Policy Evaluation?

- 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?
 - In other words, how to we take a weighted average without knowing the weights?



28

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

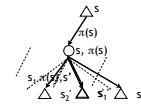
$$sample_1 = R(s, \pi(s), s'_1) + \gamma V_k^\pi(s'_1)$$

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

$$\dots$$

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

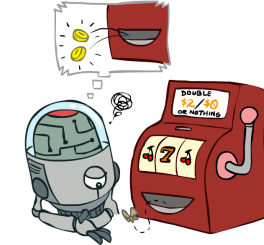
$$V_{k+1}^\pi(s) \leftarrow \frac{1}{n} \sum_i sample_i$$



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

29

Temporal Difference Learning

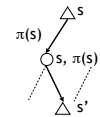


30

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!
 - Move values toward value of whatever successor occurs: running average



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

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

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

31

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)

- Decreasing learning rate (alpha) can give converging averages

32

Example: Temporal Difference Learning

States

	A	
B	C	D
	E	

Assume: $\gamma = 1$, $\alpha = 1/2$

Observed Transitions

B, east, C, -2

C, east, D, -2

	0	
0	0	8
	0	

	0	
-1	0	8
	0	

	0	
-1	3	8
	0	

$$V^\pi(s) \leftarrow (1 - \alpha) V^\pi(s) + \alpha [R(s, \pi(s), s') + \gamma V^\pi(s')]$$

33

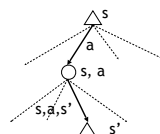
Problems with TD Value Learning

- TD value learning is a model-free way to do policy evaluation, mimicking Bellman updates with running sample averages
- However, if we want to turn values into a (new) policy, we're sunk:

$$\pi(s) = \arg \max_a Q(s, a)$$

$$Q(s, a) = \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V(s')]$$

- Idea: learn Q-values, not values
- Makes action selection model-free too!



34