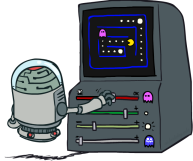


CS 188: Artificial Intelligence Reinforcement Learning II



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

The Story So Far: MDPs and RL

Known MDP: Offline Solution	
Goal	Technique
Compute V^* , Q^* , π^*	Value / policy iteration
Evaluate a fixed policy π	Policy evaluation
Unknown MDP: Model-Based	
Goal	Technique
Compute V^* , Q^* , π^*	VI/PI on approx. MDP
Evaluate a fixed policy π	PE on approx. MDP
Unknown MDP: Model-Free	
Goal	Technique
Compute V^* , Q^* , π^*	Q-learning
Evaluate a fixed policy π	Value Learning

2

Example: Direct Evaluation

Input Policy π	Observed Episodes (Training)	Output Values
	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	

Assume: $\gamma = 1$

3

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	B	C	+4
		E	-2

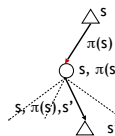
If B and E both go to C under this policy, how can their values be different?

4

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?



5

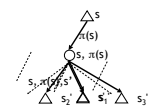
Sample-Based Policy Evaluation?

- 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

$$\begin{aligned} \text{sample}_1 &= R(s, \pi(s), s'_1) + \gamma V_k^\pi(s'_1) \\ \text{sample}_2 &= R(s, \pi(s), s'_2) + \gamma V_k^\pi(s'_2) \\ &\vdots \\ \text{sample}_n &= R(s, \pi(s), s'_n) + \gamma V_k^\pi(s'_n) \end{aligned}$$

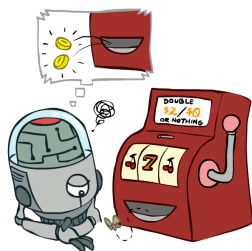
$$V_{k+1}^\pi(s) \leftarrow \frac{1}{n} \sum_i \text{sample}_i$$



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

6

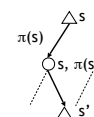
Temporal Difference Learning



7

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



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

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

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

8

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
 - Forgets about the past (distant past values were wrong anyway)
- Decreasing learning rate (α) can give converging averages

9

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	

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

10

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!

11

Active Reinforcement Learning

12

Active Reinforcement Learning

- Full reinforcement learning: optimal policies (like value iteration)
 - You don't know the transitions $T(s, a, s')$
 - You don't know the rewards $R(s, a, s')$
 - You choose the actions now
 - Goal: learn the optimal policy / values
- In this case:
 - Learner makes choices!
 - Fundamental tradeoff: exploration vs. exploitation
 - This is NOT offline planning! You actually take actions in the world and find out what happens...

13

Detour: Q-Value Iteration

- Value iteration: find successive (depth-limited) values
 - Start with $V_0(s) = 0$, which we know is right
 - Given V_k , calculate the depth $k+1$ values for all states:

$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_k(s')]$$
- But Q-values are more useful, so compute them instead
 - Start with $Q_0(s, a) = 0$, which we know is right
 - Given Q_k , calculate the depth $k+1$ q-values for all q-states:

$$Q_{k+1}(s, a) \leftarrow \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma \max_{a'} Q_k(s', a')]$$

14

Q-Learning

- Q-Learning: sample-based Q-value iteration

$$Q_{k+1}(s, a) \leftarrow \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma \max_{a'} Q_k(s', a')]$$
- Learn $Q(s, a)$ values as you go
 - Receive a sample (s, a, s', r)
 - Consider your old estimate: $Q(s, a)$
 - Consider your new sample estimate:

$$\text{sample} = R(s, a, s') + \gamma \max_{a'} Q(s', a')$$
 - Incorporate the new estimate into a running average:

$$Q(s, a) \leftarrow (1 - \alpha)Q(s, a) + (\alpha) [\text{sample}]$$

[Demo: Q-learning - gridworld (L10D2)]
[Demo: Q-learning - crawler (L10D3)]

15

Video of Demo Q-Learning -- Gridworld

16

Video of Demo Q-Learning -- Crawler

17

Q-Learning Properties

- Amazing result: Q-learning converges to optimal policy -- even if you're acting suboptimally!
- This is called off-policy learning
- Caveats:
 - You have to explore enough
 - You have to eventually make the learning rate small enough
 - ... but not decrease it too quickly
 - Basically, in the limit, it doesn't matter how you select actions (!)

[Demo: Q-learning - auto - cliff grid (L11D1)]

18

Video of Demo Q-Learning Auto Cliff Grid

