

# CS 188: Artificial Intelligence Fall 2010

## Lecture 11: Reinforcement Learning 9/30/2010

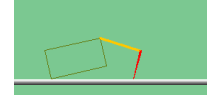
Dan Klein – UC Berkeley

Many slides over the course adapted from either Stuart Russell or Andrew Moore

## Reinforcement Learning

### Reinforcement learning:

- Still assume an 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  $\pi(s)$



[DEMO]

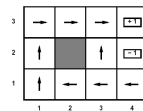
- New twist: **don't know T or R**
  - I.e. don't know which states are good or what the actions do
  - Must actually try actions and states out to learn

2

## Passive Learning

### Simplified task

- You don't know the transitions  $T(s,a,s')$
- You don't know the rewards  $R(s,a,s')$
- You are given a policy  $\pi(s)$
- Goal: learn the state values**
- ... what policy evaluation did



### In this case:

- Learner "along for the ride"
- No choice about what actions to take
- Just execute the policy and learn from experience
- We'll get to the active case soon
- This is NOT offline planning! You actually take actions in the world and see what happens...

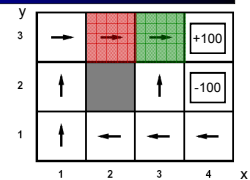
3

## Example: Direct Evaluation

[DEMO – Optimal Policy]

### Episodes:

- |                 |                 |
|-----------------|-----------------|
| (1,1) up -1     | (1,1) up -1     |
| (1,2) up -1     | (1,2) up -1     |
| (1,2) up -1     | (1,3) right -1  |
| (1,3) right -1  | (2,3) right -1  |
| (2,3) right -1  | (3,3) right -1  |
| (3,3) right -1  | (3,2) up -1     |
| (3,2) up -1     | (4,2) exit -100 |
| (3,3) right -1  | (done)          |
| (4,3) exit +100 |                 |
| (done)          |                 |



$\gamma = 1, R = -1$

$$V(2,3) \sim (96 + -103) / 2 = -3.5$$

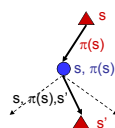
$$V(3,3) \sim (99 + 97 + -102) / 3 = 31.3$$

4

## Recap: Model-Based Policy Evaluation

### Simplified Bellman updates to calculate V for a fixed policy:

- New V is expected one-step-look-ahead using current V
- Unfortunately, need T and R



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

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

5

## Model-Based Learning

### Idea:

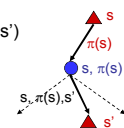
- Learn the model empirically through experience
- Solve for values as if the learned model were correct

### Simple empirical model learning

- Count outcomes for each  $s,a$
- Normalize to give estimate of  $T(s,a,s')$
- Discover  $R(s,a,s')$  when we experience  $(s,a,s')$

### Solving the MDP with the learned model

- Iterative policy evaluation, for example



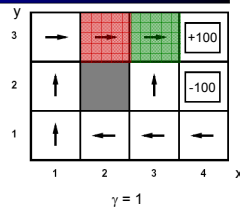
$$V_{i+1}^\pi(s) \leftarrow \sum_{s'} T(s, \pi(s), s') [R(s, \pi(s), s') + \gamma V_i^\pi(s')]$$

6

## Example: Model-Based Learning

### Episodes:

(1,1) up -1 (1,1) up -1  
 (1,2) up -1 (1,2) up -1  
 (1,2) up -1 (1,3) right -1  
 (1,3) right -1 (2,3) right -1  
 (2,3) right -1 (3,3) right -1  
 (3,3) right -1 (3,2) up -1  
 (3,2) up -1 (4,2) exit -100  
 (3,3) right -1 (done)  
 (4,3) exit +100  
 (done)



$T(<3,3>, \text{right}, <4,3>) = 1 / 3$

$T(<2,3>, \text{right}, <3,3>) = 2 / 2$

7

## Model-Free Learning

- Want to compute an expectation weighted by  $P(x)$ :

$$E[f(x)] = \sum_x P(x)f(x)$$

- Model-based: estimate  $P(x)$  from samples, compute expectation

$$x_i \sim P(x)$$

$$\hat{P}(x) = \text{count}(x)/k$$

$$E[f(x)] \approx \sum_x \hat{P}(x)f(x)$$

- Model-free: estimate expectation directly from samples

$$x_i \sim P(x)$$

$$E[f(x)] \approx \frac{1}{k} \sum_i f(x_i)$$

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

8

## Sample-Based Policy Evaluation?

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

- Who needs T and R? Approximate the expectation with samples (drawn from  $T$ !)

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

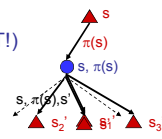
$$\text{sample}_2 = R(s, \pi(s), s'_2) + \gamma V_i^\pi(s'_2)$$

$$\dots$$

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

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

Almost! But we only actually make progress when we move to  $i+1$ .



9

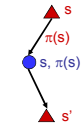
## Temporal-Difference Learning

- Big idea: learn from every experience!

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

- Temporal difference learning

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



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

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

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

10

## Exponential Moving Average

- Exponential moving average

- 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)
- Easy to compute from the running average

$$\bar{x}_n = (1 - \alpha) \cdot \bar{x}_{n-1} + \alpha \cdot x_n$$

- Decreasing learning rate can give converging averages

11

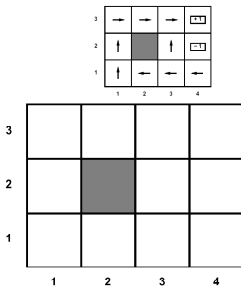
[DEMO - Grid V's]

## Example: TD Policy Evaluation

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

(1,1) up -1 (1,1) up -1  
 (1,2) up -1 (1,2) up -1  
 (1,2) up -1 (1,3) right -1  
 (1,3) right -1 (2,3) right -1  
 (2,3) right -1 (3,3) right -1  
 (3,3) right -1 (3,2) up -1  
 (3,2) up -1 (4,2) exit -100  
 (3,3) right -1 (done)  
 (4,3) exit +100  
 (done)

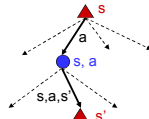
Take  $\gamma = 1, \alpha = 0.5$



12

## Problems with TD Value Learning

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



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

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

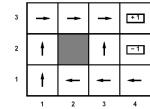
- Idea: learn Q-values directly
- Makes action selection model-free too!

13

## Active Learning

### Full reinforcement learning

- You don't know the transitions  $T(s, a, s')$
- You don't know the rewards  $R(s, a, s')$
- You can choose any actions you like
- Goal: learn the optimal policy**
- ... what value iteration did!



### 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...

14

## Detour: Q-Value Iteration

- Value iteration: find successive approx optimal values
  - Start with  $V_0(s) = 0$ , which we know is right (why?)
  - Given  $V_i$ , calculate the values for all states for depth  $i+1$ :

$$V_{i+1}(s) \leftarrow \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_i(s')]$$

- But Q-values are more useful!

- Start with  $Q_0(s, a) = 0$ , which we know is right (why?)
- Given  $Q_i$ , calculate the q-values for all q-states for depth  $i+1$ :

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

15

[DEMO – Grid Q's]

## Q-Learning

- Q-Learning: sample-based Q-value iteration
- Learn  $Q^*(s, a)$  values

- Receive a sample  $(s, a, s', r)$
- Consider your old estimate:  $Q(s, a)$
- Consider your new sample estimate:

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

$$sample = R(s, a, s') + \gamma \max_{a'} Q(s', a')$$

- Incorporate the new estimate into a running average:

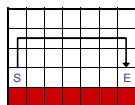
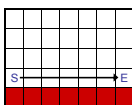
$$Q(s, a) \leftarrow (1 - \alpha)Q(s, a) + (\alpha) [sample]$$

16

[DEMO – Grid Q's]

## Q-Learning Properties

- Amazing result: Q-learning converges to optimal policy
  - If you explore enough
  - If you make the learning rate small enough
  - ... but not decrease it too quickly!
  - Basically doesn't matter how you select actions (!)
- Neat property: off-policy learning
  - learn optimal policy without following it (some caveats)



17

## Exploration / Exploitation

### Several schemes for forcing exploration

- Simplest: random actions ( $\epsilon$  greedy)
  - Every time step, flip a coin
  - With probability  $\epsilon$ , act randomly
  - With probability  $1-\epsilon$ , act according to current policy

### Problems with random actions?

- You do explore the space, but keep thrashing around once learning is done
- One solution: lower  $\epsilon$  over time
- Another solution: exploration functions

18

## Exploration Functions

- $$Q_{i+1}(s, a) \leftarrow_{\alpha} R(s, a, s') + \gamma \max_{a'} Q_i(s', a')$$
- $$Q_{i+1}(s, a) \leftarrow_{\alpha} R(s, a, s') + \gamma \max_{a'} f(Q_i(s', a'), N(s', a'))$$

## Q-Learning

- 
- |      |      |      |       |       |       |
|------|------|------|-------|-------|-------|
| 0.59 | 0.66 | 0.64 | 1.00  | 0.58  | 0.58  |
| 0.58 | 0.70 | 0.60 | 0.78  | 0.68  | 0.88  |
| 0.54 | 0.61 | 0.66 | 0.58  | 0.52  | 0.52  |
| 0.52 | 0.52 | 0.42 | -0.51 | -0.43 | -1.00 |
| 0.49 | 0.54 | 0.40 | -0.15 | 0.31  | -0.89 |
| 0.45 | 0.43 | 0.48 | 0.33  | 0.38  | -0.38 |
| 0.47 | 0.42 | 0.27 | 0.28  | 0.28  | 0.31  |
- Q-VALUES AFTER 1000 EPISODES

4