

# CS 188: Artificial Intelligence Spring 2010

## Lecture 12: Reinforcement Learning II 2/25/2010

Pieter Abbeel – UC Berkeley  
Many slides over the course adapted from either Dan Klein,  
Stuart Russell or Andrew Moore

## Announcements

- W3 Utilities: due tonight
- P3 Reinforcement Learning (RL):
  - Out tonight, due Thursday next week
  - You will get to apply RL to:
    - Gridworld agent
    - Crawler
    - Pac-man

## 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  $\pi(s)$   $S \rightarrow A$
- 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

## The Story So Far: MDPs and RL

Things we know how to do:	Techniques:
<ul style="list-style-type: none"> <li>▪ If we know the MDP                             <ul style="list-style-type: none"> <li>▪ Compute <math>V^*, Q^*, \pi^*</math> exactly</li> <li>▪ Evaluate a fixed policy <math>\pi</math></li> </ul> </li> </ul>	<ul style="list-style-type: none"> <li>▪ Model-based DPs                             <ul style="list-style-type: none"> <li>▪ Value and policy iteration</li> <li>▪ Policy evaluation</li> </ul> </li> </ul>
<ul style="list-style-type: none"> <li>▪ If we don't know the MDP                             <ul style="list-style-type: none"> <li>▪ We can estimate the MDP then solve</li> <li>▪ We can estimate <math>V</math> for a fixed policy <math>\pi</math></li> <li>▪ We can estimate <math>Q^*(s,a)</math> for the optimal policy while executing an exploration policy</li> </ul> </li> </ul>	<ul style="list-style-type: none"> <li>▪ Model-based RL</li> <li>▪ Model-free RL:                             <ul style="list-style-type: none"> <li>▪ Value learning</li> <li>▪ Q-learning</li> </ul> </li> </ul>

*arg max<sub>a</sub> Q\*(s,a)*

## 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!

## 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...

3	←	←	←	☐
2	↑	■	↑	☐
1	↑	←	←	←
	1	2	3	4

### Detour: Q-Value Iteration

*has access to T, R*

- 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  $i+1$ :

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

*observed (s,a,s') empirical update  $\frac{R(s,a,s')}{\sum_{s'} T(s,a,s')} = \frac{R(s,a,s')}{\gamma \max_{a'} Q_i(s',a')}$*

- 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  $i+1$ :

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

### Q-Learning

*sample estimate:  $x_i \sim P(x_i)$*

$$E[f(x)] = \sum P(x_i) f(x_i)$$

$$E[f(x)] = \frac{1}{n} \sum f(x_i)$$

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

### Q-Learning Properties

- Amazing result: Q-learning converges to optimal policy
  - If you explore enough  $\sim$  all  $(s,a)$  visited  $\infty$  often
  - If you make the learning rate small enough  $\frac{1}{n}$
  - ... but not decrease it too quickly!  $\sum \frac{1}{n} = \infty$
  - Basically doesn't matter how you select actions (I)
- Neat property: off-policy learning
  - learns optimal Q-values, not the values of the policy you are following

### 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 e.g.  $\max_{a'} Q(s,a)$
- Regret: expected gap between rewards during learning and rewards from optimal action
  - Q-learning with random actions will converge to optimal values, but possibly very slowly, and will get low rewards on the way
  - Results will be optimal but regret will be large
  - How to make regret small?

### Exploration Functions

- When to explore
  - Random actions: explore a fixed amount
  - Better ideas: explore areas whose badness is not established, explore less over time
- One way: exploration function
  - Takes a value estimate and a count, and returns an optimistic utility, e.g.  $f(u, n) = u + \frac{k}{n}$  (exact form not important)

$$Q_{i+1}(s, a) \leftarrow \alpha R(s, a, s') + \gamma \max_{a'} Q_i(s', a')$$

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

### Q-Learning

- Q-learning produces tables of q-values:

0.59	0.66	0.64	1.00
0.58	0.70	0.60	0.88
0.54	0.66	0.58	0.58
0.61		-0.51	-1.00
0.52	0.52	0.42	-0.43
0.49		0.35	-0.89
0.54	0.40	0.31	-0.89
0.45	0.43	0.39	-0.38
0.47	0.42	0.27	0.11

Q-VALUES AFTER 1000 EPISODES



