







RL Framework

- Learn by "trial and error"
- No assumptions about model
- No assumptions about reward function
- Assumes:
 - True state is known at all times
 - Immediate reward is known
 - Discount is known



- Problem: We don't know probability of answering correctly
- Solution:
 - Buy the home version of the game



- Practice on the home game to refine our strategy
- Deploy strategy when we play the real game

Model Learning Approach

- Learn model, solve
- How to learn a model:
 - Take action a in state s, observe s'
 - Take action a in state s, n times
 - Observe s' m times
 - P(s'|s,a) = m/n
 - Fill in transition matrix for each action
 - Compute avg. reward for each state
- Solve learned model as an MDP (previous lecture)



















Learning Rates in RL in Practice Maintain a per-state count N[s] Learning rate is function of N[s], α(N[s]) Sufficient to satisfy theory: α(N[s])=1/N(s) Often viewed as too slow α drops quickly

- Convergence is slow
- In practice, often a floor on, α , e.g., α = 0.01
- Floor leads to faster learning, but less stability























Using TD for Control

• Recall value iteration:

$$V^{i+1}(s) = \max_{a} R(s,a) + \gamma \sum_{s'} P(s'|s,a) V^{i}(s')$$

• Why not pick the maximizing **a** and then do:

$$V(s) = V(s) + \alpha(N(s))(r + \gamma V(s') - V(s))$$

- s' is the observed next state after taking action a







Q-Learning Demo (see video from Zoom)

- Keyboard controlled Q-learning agent using robot grid world from R&N (and previous MDP slides)
- Differences:
 - Discount of 0.9
 - No step cost (previously -0.04)
 - Bad state has value 0 (previously -1)
- How to think of terminal states:
 - Make transition to another state (absorbing state) with value that is always 0
 - Problem then resets to start (no transition from absorbing state to start)



Brief Comments on Learning from Demonstration

- LfD is a powerful method to convey human expertise to (ro)bots
- Useful for imitating human policies
- Less useful for surpassing human ability (but can smooth out noise in human demos)
- Used, e.g., for acrobatic helicopter flight



Exploration vs. Exploitation

- Greedy strategy purely **exploits** its current knowledge
 - The quality of this knowledge improves only for those states that the agent observes often
- A good learner must perform exploration in order to improve its knowledge about states that are not often observed
 - But pure exploration is useless (and costly) if it is never exploited



Exploration vs. Exploitation in Theory and Practice

- Can assign an "exploration bonus" to parts of the world (or state-action combinations) you haven't experienced much
 - Versions of this are provably efficient, e.g., R-Max (will eventually learn the optimal policy requiring polynomial effort in size of problem)
 - Works for small state spaces
- In practice ε-greedy action selection is used most often
 - Choose greedy action w.p. 1- ϵ
 - Choose random action w.p. ϵ







Linear Regression Review

• Define a set of basis functions (vectors)

 $\varphi_1(s), \varphi_2(s)...\varphi_k(s)$

• Approximate f with a weighted combination of these

$$g(s;w) = \sum_{j=1}^{\kappa} w_j \varphi_j(s)$$

• Example: Space of quadratic functions:

$$\varphi_1(s) = 1, \varphi_2(s) = s, \varphi_3(s) = s^2$$

• Orthogonal projection minimizes SSE



Properties of approximate RL

- Exact case (tabular representation) = special case
- Can be combined with Q-learning
- Convergence not guaranteed
 - Policy evaluation with linear function approximation converges if samples are drawn "on policy"
 - In general, convergence is not guaranteed
 - Chasing a moving target
 - Errors can compound
- Success has traditionally required very carefully chosen features
- Deepmind has recently had success using no feature engineering but lots of training data



Conclusions

- Reinforcement learning solves an MDP
- Converges for exact value function representation
- Can be combined with approximation methods
- Good results require good features and/or lots of data