

CS 188: Artificial Intelligence Fall 2010

Lecture 9: MDPs 9/23/2010

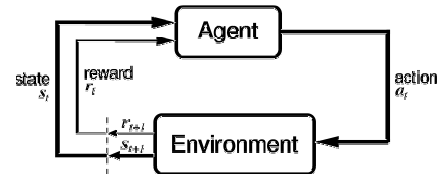
Dan Klein – UC Berkeley

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

Reinforcement Learning

[DEMOS]

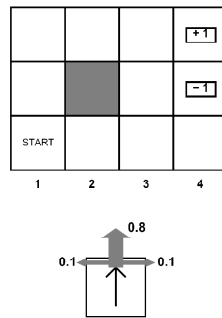
- Basic idea:
 - Receive feedback in the form of **rewards**
 - Agent's utility is defined by the reward function
 - Must (learn to) act so as to **maximize expected rewards**



Grid World

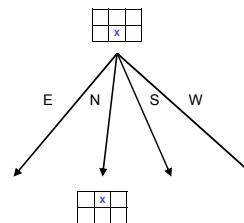
[DEMO – Gridworld Intro]

- The agent lives in a grid
- Walls block the agent's path
- The agent's actions do not always go as planned:
 - 80% of the time, the action North takes the agent North (if there is no wall there)
 - 10% of the time, North takes the agent West; 10% East
 - If there is a wall in the direction the agent would have been taken, the agent stays put
- Small "living" reward each step
- Big rewards come at the end
- Goal: maximize sum of rewards*

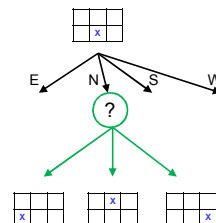


Grid Futures

Deterministic Grid World

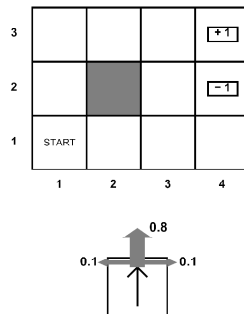


Stochastic Grid World



Markov Decision Processes

- An MDP is defined by:
 - A set of states $s \in S$
 - A set of actions $a \in A$
 - A transition function $T(s,a,s')$
 - Prob that a from s leads to s'
 - i.e., $P(s' | s,a)$
 - Also called the model
 - A reward function $R(s, a, s')$
 - Sometimes just $R(s)$ or $R(s')$
 - A start state (or distribution)
 - Maybe a terminal state
- MDPs are a family of non-deterministic search problems
 - Reinforcement learning: MDPs where we don't know the transition or reward functions



What is Markov about MDPs?

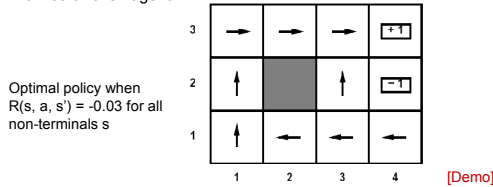
- Andrey Markov (1856-1922)
- "Markov" generally means that given the present state, the future and the past are independent
- For Markov decision processes, "Markov" means:



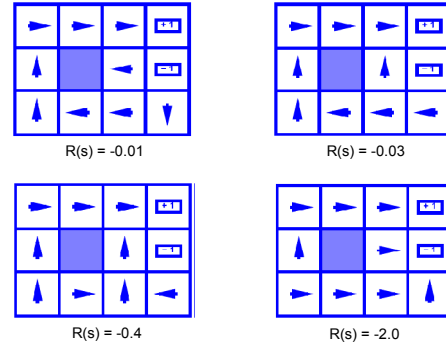
$$P(S_{t+1} = s' | S_t = s_t, A_t = a_t, S_{t-1} = s_{t-1}, A_{t-1}, \dots, S_0 = s_0) \\ = \\ P(S_{t+1} = s' | S_t = s_t, A_t = a_t)$$

Solving MDPs

- In deterministic single-agent search problems, want an optimal **plan**, or sequence of actions, from start to a goal
- In an MDP, we want an optimal **policy** $\pi^*: S \rightarrow A$
 - A policy π gives an action for each state
 - An optimal policy maximizes expected utility if followed
 - Defines a reflex agent

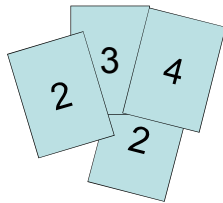


Example Optimal Policies



Example: High-Low

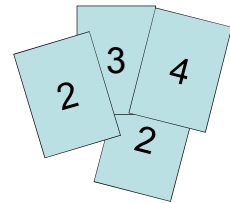
- Three card types: 2, 3, 4
- Infinite deck, twice as many 2's
- Start with 3 showing
- After each card, you say "high" or "low"
- New card is flipped
- If you're right, you win the points shown on the new card
- Ties are no-ops
- If you're wrong, game ends
- Why not use expectimax?
 - #1: get rewards as you go
 - #2: you might play forever!



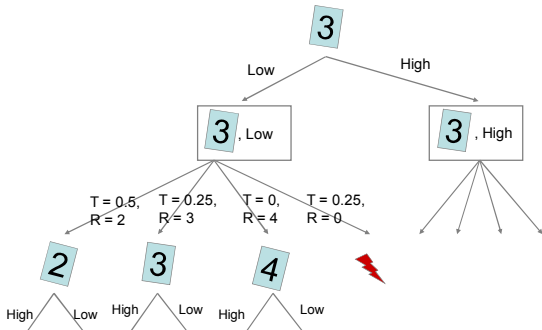
9

High-Low as an MDP

- States: 2, 3, 4, done
- Actions: High, Low
- Model: $T(s, a, s')$:
 - $P(s'=4 | 4, \text{Low}) = 1/4$
 - $P(s'=3 | 4, \text{Low}) = 1/4$
 - $P(s'=2 | 4, \text{Low}) = 1/2$
 - $P(s'=\text{done} | 4, \text{Low}) = 0$
 - $P(s'=4 | 4, \text{High}) = 1/4$
 - $P(s'=3 | 4, \text{High}) = 0$
 - $P(s'=2 | 4, \text{High}) = 0$
 - $P(s'=\text{done} | 4, \text{High}) = 3/4$
 - ...
- Rewards: $R(s, a, s')$:
 - Number shown on s' if $s \neq s'$
 - 0 otherwise
- Start: 3



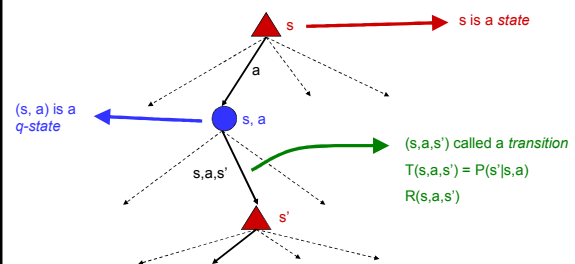
Example: High-Low



11

MDP Search Trees

- Each MDP state gives an expectimax-like search tree



12

Utilities of Sequences

- In order to formalize optimality of a policy, need to understand utilities of sequences of rewards

- Typically consider **stationary preferences**:

$$\begin{aligned} [r, r_0, r_1, r_2, \dots] &> [r, r'_0, r'_1, r'_2, \dots] \\ &\Leftrightarrow \\ [r_0, r_1, r_2, \dots] &> [r'_0, r'_1, r'_2, \dots] \end{aligned}$$

- Theorem: only two ways to define stationary utilities**

- Additive utility:

$$U([r_0, r_1, r_2, \dots]) = r_0 + r_1 + r_2 + \dots$$

- Discounted utility:

$$U([r_0, r_1, r_2, \dots]) = r_0 + \gamma r_1 + \gamma^2 r_2 + \dots$$

13

Infinite Utilities?!

- Problem: infinite state sequences have infinite rewards

- Solutions:**

- Finite horizon:

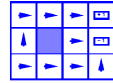
- Terminate episodes after a fixed T steps (e.g. life)
- Gives nonstationary policies (π depends on time left)

- Absorbing state: guarantee that for every policy, a terminal state will eventually be reached (like "done" for High-Low)

- Discounting: for $0 < \gamma < 1$

$$U([r_0, \dots, r_\infty]) = \sum_{t=0}^{\infty} \gamma^t r_t \leq R_{\max} / (1 - \gamma)$$

- Smaller γ means smaller "horizon" – shorter term focus



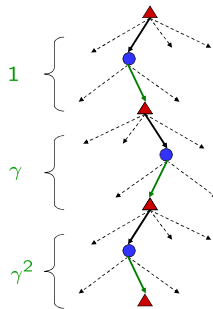
14

Discounting

- Typically discount rewards by $\gamma < 1$ each time step

- Sooner rewards have higher utility than later rewards

- Also helps the algorithms converge



15

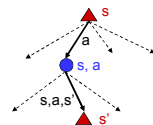
Recap: Defining MDPs

- Markov decision processes:**

- States S
- Start state s_0
- Actions A
- Transitions $P(s'|s,a)$ (or $T(s,a,s')$)
- Rewards $R(s,a,s')$ (and discount γ)

- MDP quantities so far:**

- Policy = Choice of action for each state
- Utility (or return) = sum of discounted rewards



16

Optimal Utilities

- Fundamental operation: compute the values (optimal expectimax utilities) of states s

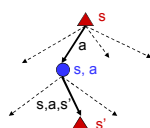
- Why? Optimal values define optimal policies!

- Define the value of a state s :
 $V(s)$ = expected utility starting in s and acting optimally

- Define the value of a q-state (s,a) :
 $Q(s,a)$ = expected utility starting in s , taking action a and thereafter acting optimally

- Define the optimal policy:
 $\pi^*(s)$ = optimal action from state s

[DEMO – Grid Values]



3	0.512	0.868	0.912	1.0	3
2	0.762	0.660	0.660	1.0	2
1	0.705	0.655	0.611	0.388	1
	1	2	3	4	

17

The Bellman Equations

- Definition of "optimal utility" leads to a simple one-step lookahead relationship amongst optimal utility values:

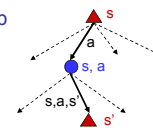
Optimal rewards = maximize over first action and then follow optimal policy

- Formally:

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

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

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



18

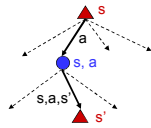
Solving MDPs

- We want to find the **optimal policy** π^*
- Proposal 1: modified expectimax search, starting from each state s :

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

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

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



19

Why Not Search Trees?

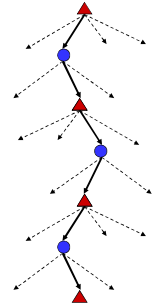
- Why not solve with expectimax?

Problems:

- This tree is usually infinite (why?)
- Same states appear over and over (why?)
- We would search once per state (why?)

Idea: Value iteration

- Compute optimal values for all states all at once using successive approximations
- Will be a bottom-up dynamic program similar in cost to memoization
- Do all planning offline, no replanning needed!

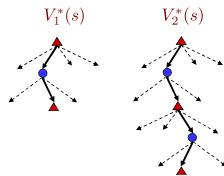


20

Value Estimates

- Calculate estimates $V_k^*(s)$

- Not the optimal value of s !
- The optimal value considering only next k time steps (k rewards)
- As $k \rightarrow \infty$, it approaches the optimal value



- Almost solution: recursion (i.e. expectimax)
- Correct solution: dynamic programming

[DEMO -- V_k]

21

Value Iteration

Idea:

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

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

- This is called a **value update** or **Bellman update**
- Repeat until convergence

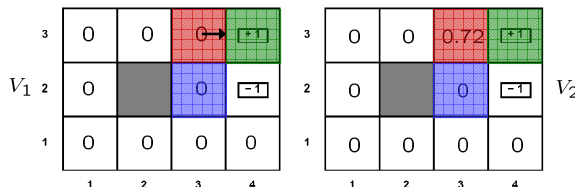
Theorem: will converge to unique optimal values

- Basic idea: approximations get refined towards optimal values
- Policy may converge long before values do

22

Example: Bellman Updates

Example: $\gamma=0.9$, living reward=0, noise=0.2



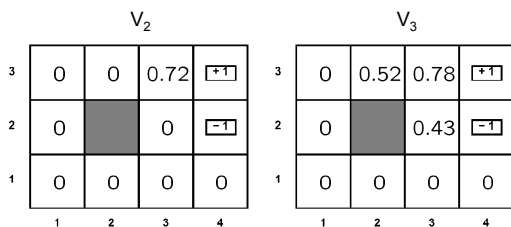
$$V_{i+1}(s) = \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_i(s')]$$

$$V_2((3,3)) = \sum_{s'} T((3,3), \text{right}, s') [R((3,3), \text{right}, s') + 0.9 V_1(s')] = 0.9 [0.8 \cdot 1 + 0.1 \cdot 0 + 0.1 \cdot 0]$$

max happens for $a=\text{right}$, other actions not shown

23

Example: Value Iteration



- Information propagates outward from terminal states and eventually all states have correct value estimates

[DEMO]

24

Convergence*

- Define the max-norm: $\|U\| = \max_s |U(s)|$

- Theorem: For any two approximations U and V

$$\|U^{t+1} - V^{t+1}\| \leq \gamma \|U^t - V^t\|$$

- I.e. any distinct approximations must get closer to each other, so, in particular, any approximation must get closer to the true U and value iteration converges to a unique, stable, optimal solution

- Theorem:

$$\|U^{t+1} - U^t\| < \epsilon, \Rightarrow \|U^{t+1} - U\| < 2\epsilon\gamma/(1 - \gamma)$$

- I.e. once the change in our approximation is small, it must also be close to correct

25