CS 343: Artificial Intelligence

Markov Decision Processes



Profs. Peter Stone and Yuke Zhu, The University of Texas at Austin

[These slides based on those of Dan Klein and Pieter Abbeel for CS188 Intro to AI at UC Berkeley. All CS188 materials are available at http://ai.berkeley.edu.]

Announcements

- Homework 3: Games
 - Due 3/1 at 11:59 pm.
- Programming Project 2: Multi-Agent Pacman
 - Due 3/3 at 11:59 pm.
- Homework 4: MDPs
 - Due 3/8 at 11:59 pm.
- Assignment Grading
 - To be released by the end of this week
- Midterm Exam
 - Details will follow

Example: Grid World

- A maze-like problem
 - The agent lives in a grid
 - Walls block the agent's path
- Noisy movement: actions do not always go as planned
 - 80% of the time, the action has the intended effect (if there is no wall there)
 - 20% of the time an adjacent action occurs instead. Ex: North has 10% chance of East and 10% chance of West
 - If there is a wall in the direction the agent would have been taken, the agent stays put
- The agent receives rewards each time step
 - Small "living" reward each step (can be negative)
 - Big rewards come at the end (good or bad)
- Goal: maximize sum of rewards



Grid World Actions

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')
 - Probability that a from s leads to s', i.e., P(s' | s, a)
 - Also called the model or the dynamics
 - A reward function R(s, a, s')
 - Sometimes just R(s) or R(s')
 - A start state
 - Maybe a terminal state
- MDPs are non-deterministic search problems
 - One way to solve them is with expectimax search
 - ...but with modification to allow rewards along the way
 - We'll have a new, more efficient tool soon



What is Markov about MDPs?

- "Markov" generally means that given the present state, the future and the past are independent
- For Markov decision processes, "Markov" means action outcomes depend only on the current state

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

$$P(S_{t+1} = s' | S_t = s_t, A_t = a_t)$$

=

 This is just like search, where the successor function could only depend on the current state (not the history)



Andrey Markov (1856-1922)

Policies

- In deterministic single-agent search problems, we wanted an optimal plan, or sequence of actions, from start to a goal
- For MDPs, we want an optimal policy $\pi^*: S \rightarrow A$
 - A policy π gives an action for each state
 - An optimal policy is one that maximizes expected utility if followed
 - An explicit policy defines a reflex agent
- Expectimax didn't compute entire policies
 - It computed the action for a single state only



Optimal policy when R(s, a, s') = -0.03 for all non-terminals s

Example: Racing

- A robot car wants to travel far, quickly
- Three states: Cool, Warm, Overheated
- Two actions: *Slow*, *Fast*



Racing Search Tree



MDP Search Trees



Utilities of Sequences

- What preferences should an agent have over reward sequences?
- More or less? [1, 2, 2] or [2, 3, 4]
- Now or later? [0, 0, 1] or [1, 0, 0]



Discounting

- It's reasonable to maximize the sum of rewards
- It's also reasonable to prefer rewards now to rewards later
- One solution: values of rewards decay exponentially



Discounting

- How to discount?
 - Each time we descend a level, we multiply in the discount once
- Why discount?
 - Sooner rewards probably do have higher utility than later rewards
 - Also helps our algorithms converge
- Example: discount of 0.5
 - U([1,2,3]) = 1*1 + 0.5*2 + 0.25*3
 - U([1,2,3]) < U([3,2,1])



Infinite Utilities?!

- Problem: What if the game lasts forever? Do we get infinite rewards?
- Solutions:
 - Discounting: use $0 < \gamma < 1$
 - Smaller γ means smaller "horizon" shorter term focus

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

- Finite horizon: (similar to depth-limited search)
 - 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 "overheated" for racing)



Recap: Defining MDPs

Markov decision processes:

- Set of states S
- Start state s₀
- Set of actions A
- Transitions P(s'|s,a) (or T(s,a,s'))
- Rewards R(s,a,s') (and discount γ)



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



Optimal Quantities

- The value (utility) of a state s:
 - V*(s) = expected utility starting in s and acting optimally
- The value (utility) of a q-state (s,a):

Q*(s,a) = expected utility starting out having taken action a from state s and (thereafter) acting optimally

• The optimal policy:

 $\pi^*(s)$ = optimal action from state s



Optimal Quantities

• The value (utility) of a state s:

V*(s) = expected utility starting in s and acting optimally

• The value (utility) of a q-state (s,a):

Q^{*}(s,a) = expected utility starting out having taken action a from state s and (thereafter) acting optimally

• The optimal policy:

 $\pi^*(s)$ = optimal action from state s



Values of States

- Fundamental operation: compute the (expectimax) value of a state
 - Expected utility under optimal action
 - Average sum of (discounted) rewards
 - This is just what expectimax computed!
- Recursive definition of (optimal) value:

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

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



The Bellman Equations

 Definition of "optimal utility" via expectimax recurrence gives a simple one-step lookahead relationship amongst optimal utility values

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



 These are the Bellman equations, and they characterize optimal values in a way we'll use over and over

Time-Limited Values

- Key idea: time-limited values
- Define V_k(s) to be the optimal value of s if the game ends in k more time steps
 - Equivalently, it's what a depth-k expectimax would give from s





Value Iteration

- Start with V₀(s) = 0: no time steps left means an expected reward sum of zero
- Given vector of V_k(s) values, do one step of expectimax from each state:

$$V_{k+1}(s) \leftarrow \max_{a} \sum_{s'} T(s, a, s') \left[R(s, a, s') + \gamma V_k(s') \right]$$

- Repeat until convergence
- Complexity of each iteration: O(S²A)
- Theorem: will converge to unique optimal values
 - Basic idea: approximations get refined towards optimal values
 - Policy may converge long before values do



Example: Value Iteration



Gridworld Values V*

Gridworld Display			
0.64 ▶	0.74 ▸	0.85 →	1.00
• 0.57		•	-1.00
• 0.49	∢ 0.43	• 0.48	∢ 0.28
VALUES AFTER 100 ITERATIONS			

Gridworld: Q*





Profs. Peter Stone and Yuke Zhu — The University of Texas at Austin

[These slides based on those of Dan Klein and Pieter Abbeel for CS188 Intro to AI at UC Berkeley. All CS188 materials are available at http://ai.berkeley.edu.]

Value Iteration

Bellman equations characterize the optimal values:

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

• Value iteration computes them:

$$V_{k+1}(s) \leftarrow \max_{a} \sum_{s'} T(s, a, s') \left[R(s, a, s') + \gamma V_k(s') \right]$$

- Value iteration is just a fixed point solution method
 - ... though the V_k vectors are also interpretable as time-limited values



Fixed Policies



- Expectimax trees max over all actions to compute the optimal values
- If we fixed some policy $\pi(s)$, then the tree would be simpler only one action per state
 - ... though the tree's value would depend on which policy we fixed

Utilities for a Fixed Policy

- Another basic operation: compute the utility of a state s under a fixed (generally non-optimal) policy
- Define the utility of a state s, under a fixed policy π:
 V^π(s) = expected total discounted rewards starting in s and following π
- Recursive relation (one-step look-ahead / Bellman equation):

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



Utilities for a Fixed Policy

- Another basic operation: compute the utility of a state s under a fixed (generally non-optimal) policy
- Define the utility of a state s, under a fixed policy π:
 V^π(s) = expected total discounted rewards starting in s and following π
- Recursive relation (one-step look-ahead / Bellman equation):

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

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



Example: Policy Evaluation

Always Go Right

Always Go Forward



Example: Policy Evaluation

Always Go Right

-10.00	100.00	-10.00
-10.00	1.09 🕨	-10.00
-10.00	-7.88 🕨	-10.00
-10.00	-8.69 ▶	-10.00

Always Go Forward

-10.00	100.00	-10.00
-10.00	▲ 70.20	-10.00
-10.00	▲ 48.74	-10.00
10.00	▲ 33.30	-10.00

Policy Evaluation

- How do we calculate the V's for a fixed policy π ?
- Idea 1: Turn recursive Bellman equations into updates (like value iteration)

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

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

- Efficiency: O(S²) per iteration
- Idea 2: Without the maxes, the Bellman equations are just a linear system
 - Solve with Matlab (or your favorite linear system solver)



Computing Actions from Values

- Let's imagine we have the optimal values V*(s)
- How should we act?
 - It's not obvious!
- We need to do a mini-expectimax (one step)



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

• This is called **policy extraction**, since it gets the policy implied by the values

Computing Actions from Q-Values

- Let's imagine we have the optimal q-values:
- How should we act?
 - Completely trivial to decide!

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



- Important lesson: actions are easier to select from q-values than values!
- In fact, you don't even need a model!

Problems with Value Iteration

Value iteration repeats the Bellman updates:

$$V_{k+1}(s) \leftarrow \max_{a} \sum_{s'} T(s, a, s') \left[R(s, a, s') + \gamma V_k(s') \right]$$

- Problem 1: It's slow O(S²A) per iteration
- Problem 2: The "max" at each state rarely changes
- Problem 3: The policy often converges long before the values



Policy Iteration

- Alternative approach for optimal values:
 - Step 1: Policy evaluation: calculate utilities for some fixed policy (not optimal utilities!) until convergence
 - Step 2: Policy improvement: update policy using one-step look-ahead with resulting converged (but not optimal!) utilities as future values
 - Repeat steps until policy converges
- This is policy iteration
 - It's still optimal!
 - Can converge (much) faster under some conditions

Policy Iteration

- Evaluation: For fixed current policy π , find values with policy evaluation:
 - Iterate until values converge:

$$V_{k+1}^{\pi_i}(s) \leftarrow \sum_{s'} T(s, \pi_i(s), s') \left[R(s, \pi_i(s), s') + \gamma V_k^{\pi_i}(s') \right]$$

- Improvement: For fixed values, get a better policy using policy extraction
 - One-step look-ahead:

$$\pi_{i+1}(s) = \arg\max_{a} \sum_{s'} T(s, a, s') \left[R(s, a, s') + \gamma V^{\pi_i}(s') \right]$$

0 0	Gridworl	d Display			
0.00	0.00	0.00	0.00		
^		^			
0.00		0.00	0.00		
		^	^		
0.00	0.00	0.00	0.00		
VALUI	VALUES AFTER O TTERATIONS				

00	O O Gridworld Display			
	•	• 0.00	0.00)	1.00
	• 0.00		∢ 0.00	-1.00
	•	•	• 0.00	0.00
	VALUES AFTER 1 ITERATIONS			

Gridworld Display					
•	0.00 >	0.72 →	1.00		
0.00		0.00	-1.00		
^	^	^			
0.00	0.00	0.00	0.00		
			-		
VALUE	VALUES AFTER 2 ITERATIONS				

k=3

0	0	Gridworl	d Display		
	0.00)	0.52 →	0.78)	1.00	
	• 0.00		• 0.43	-1.00	
	•	•	•	0.00	
	VALUES AFTER 3 ITERATIONS				

k=4

00		Gridworl	d Display		
	0.37 →	0.66)	0.83)	1.00	
	• 0.00		• 0.51	-1.00	
	• 0.00	0.00 →	• 0.31	∢ 0.00	
	VALUES AFTER 4 ITERATIONS				

0 0	C C Cridworld Display				
	0.51)	0.72 ▸	0.84 →	1.00	
	0. 27		• 0.55	-1.00	
			•		
	0.00	0.22 ♪	0.37	◆ 0.13	
	VALUES AFTER 5 ITERATIONS				

00	Gridworl	d Display		
0.59)	0.73 →	0.85)	1.00	
• 0.41		• 0.57	-1.00	
• 0.21	0.31 →	• 0.43	∢ 0.19	
VALUES AFTER 6 ITERATIONS				

00	Gridworl	d Display	-	
0.62)	0.74 →	0.85)	1.00	
• 0.50		• 0.57	-1.00	
▲ 0.34	0.36)	▲ 0.45	◀ 0.24	
VALUES AFTER 7 ITERATIONS				

0 0	Gridworl	d Display	
0.63)	0.74 ▸	0.85)	1.00
• 0.53		• 0.57	-1.00
• 0.42	0.39 →	▲ 0.46	∢ 0.26
VALUE	S AFTER	8 ITERA	FIONS

Gridworld Display			
0.64 →	0.74 →	0.85)	1.00
• 0.55		• 0.57	-1.00
• 0.46	0.40 →	• 0.47	∢ 0.27
VALUES AFTER 9 ITERATIONS			

○ ○ Gridworld Display			
0.64)	0.74 →	0.85)	1.00
• 0.56		• 0.57	-1.00
• 0.48	∢ 0.41	• 0.47	∢ 0.27
VALUES AFTER 10 ITERATIONS			

O O O Gridworld Display				
	0.64 →	0.74)	0.85)	1.00
	• 0.56		• 0.57	-1.00
	•	◀ 0.42	• 0.47	∢ 0.27
VALUES AFTER 11 ITERATIONS				

O O O Gridworld Display				
	0.64 ▸	0.74 ▸	0.85)	1.00
	0.57		0.57	-1.00
	• 0.49	◀ 0.42	• 0.47	∢ 0.28
	VALUES AFTER 12 ITERATIONS			

Gridworld Display			
0.64)	0.74 →	0.85)	1.00
• 0.57		• 0.57	-1.00
• 0.49	◀ 0.43	▲ 0.48	∢ 0.28
VALUES AFTER 100 ITERATIONS			

Comparison

- Both value iteration and policy iteration compute the same thing (all optimal values)
- In value iteration:
 - Every iteration updates both the values and (implicitly) the policy
 - We don't track the policy, but taking the max over actions implicitly recomputes it
- In policy iteration:
 - We do several passes that update utilities with fixed policy (each pass is fast because we consider only one action, not all of them)
 - After the policy is evaluated, a new policy is chosen (slow like a value iteration pass)
 - The new policy will be better (or we're done)
- Both are dynamic programs for solving MDPs

Summary: MDP Algorithms

- So you want to....
 - Compute optimal values: use value iteration or policy iteration
 - Compute values for a particular policy: use policy evaluation
 - Turn your values into a policy: use policy extraction (one-step lookahead)

These all look the same!

- They basically are they are all variations of Bellman updates
- They all use one-step lookahead expectimax fragments
- They differ only in whether we plug in a fixed policy or max over actions

Double Bandits



Double-Bandit MDP



Offline Planning

Solving MDPs is offline planning

- You determine all quantities through computation
- You need to know the details of the MDP
- You do not actually play the game!

No discount 100 time steps Both states have the same value



Online Planning

Rules changed! Red's win chance is different.



Next Time: Reinforcement Learning!