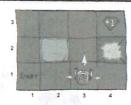
CS 188: Artificial Intelligence Markov Decision Processes II



Instructors: Dan Klein and Pieter Abbeel --- University of California, Berkeley
These sides core created to Dan Klein and Pieter Abbeel for CS188 toro to Al at U.C. Beliefer. Ad CS188 magnetic are available at http://du.berheing.adu."

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 North takes the agent North
- 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
- . The agent receives rewards each time step
- Small "living" reward each step (can be riegative)
- Big rewards come at the end (good or bad)
- Goal: maximize sum of (discounted) rewards



we assume a weighted averaging

Recap: MDPs

- Markov decision processes:
- States 5
- · Actions A
- Transitions P(s'|s,a) (or T(s,a,s'))
- Rewards R(s,a,s') (and discount y)
- Start state s_o



- Policy = map of states to actions
- Utility = sum of discounted rewards
- * Values = expected future utility from a state (max node) woler oftinal or
- Q-Values = expected future utility from a q-state (chance node)

chance node



Reward R(s,a,s') occurs during each transition a + is instant, not expectation

Gridworld: Q*

Optimal Quantities

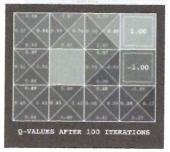
- The value (utility) of a state s:
 V*(s) = expected utility starting in s and acting optimally
- The value (utility) of a q-state (s,a):
 Q'(s,a) = expected utility starting out having taken action a from state s and (thereafter) acting optimally
- The optimal policy:
 π*(s) = optimal action from state s
- · Usually good is to find The
- · Values + Q-values are



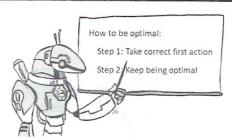
.

0.64 , 0.74 ,		
0.57	0.57	-1.00
0.49 . 0.43	0.48	< 0.28
VALUES AFTER	100 ITER	ATIONS

Gridworld Values V*



The Bellman Equations



The Bellman Equations

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

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

$$Q^{*}(s, n) \equiv \sum \Gamma(s, n, s') \cdot R(s, n, s') + \gamma V^{*}(s')$$

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

- These are the Bellman equations, and they characterize optimal values in a way we'll use over and over
- · Just a system of eq., not an algorithm by itself.

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:

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

- Value iteration is just a fixed point solution method
 ... though the V, vectors are also interpretable as time-limited values
- · Start w/ Vo(s') which always = 0

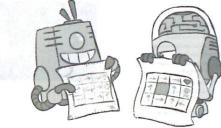
 $V_k(s)$

Policy Methods

Policy Evaluation

- · How do we know the V, vectors are going to converge?
- Case 1: If the tree has maximum depth M, then V_{st} holds the actual untruncated values
- . Case 2. If the discount is less than 1.
 - \bullet . Sketch: For any state V_k and V_{k-1} can be viewed as depth k+1 expectimax results in nearly identical search trees
 - . The difference is that on the bottom layer, V_{tri} has actual rewards while V, has zeros
 - . That last layer is at best all Russ
 - . It is at worst R
 - But everything is discounted by y^a that far out
 - So V_i and V_{i+1} are at most yⁱ max |R| different
 - · So as k increases, the values converge







- · Given a policy, how well · Specifically, for a given policy, finding V(s) for each s
 - Example: Policy Evaluation

Fixed Policies

Utilities for a Fixed Policy



* Expectimax trees max over all actions to compute the optimal values

... though the tree's value would depend on which policy we fixed







- 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''(s) = expected total discounted rewards starting in s and following π
- · Recursive relation (one-step look-ahead / Bellman equation):

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

according to

Always Go Right



Always Go Forward



Example: Policy Evaluation

If we fixed some policy n(s), then the tree would be simpler - only one action per state

Policy Evaluation

Policy Extraction

Always Go Right



Always Go Forward



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

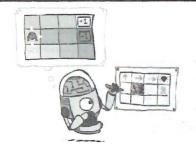
$$V_0^{\pi}(s) = 0$$

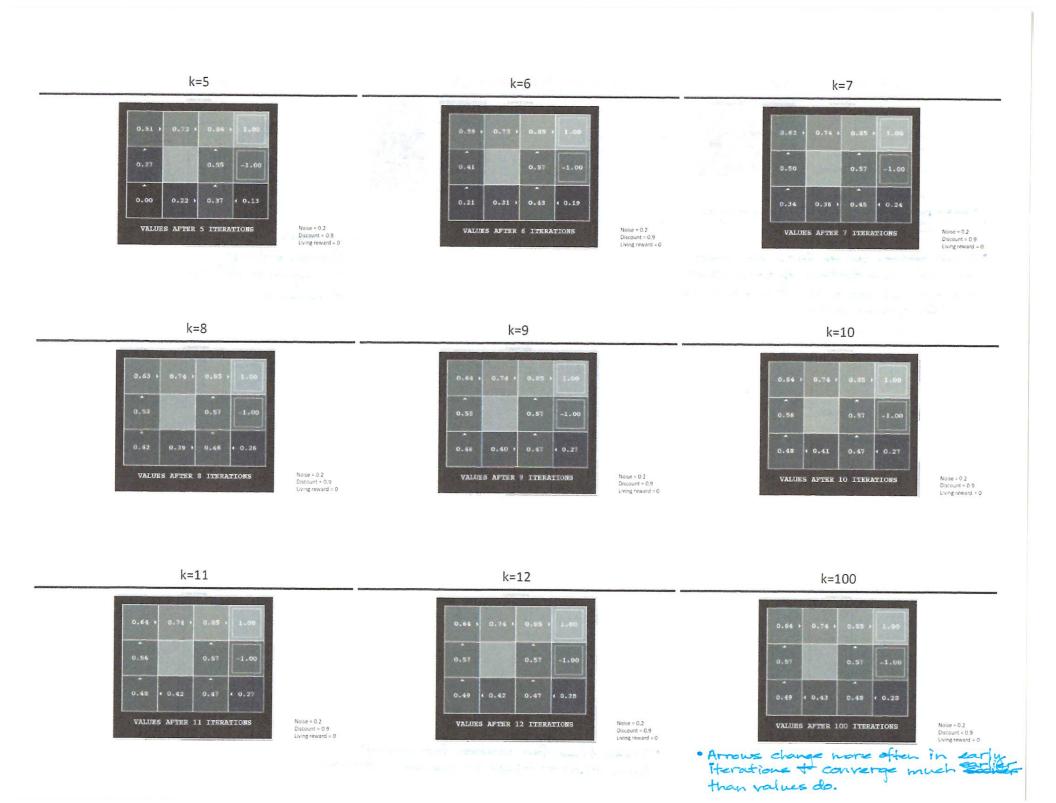


- Efficiency: O(S²) per iteration
 val iter: O(as²)

 $\pi(s)$

- Idea 2: Without the maxes, the Bellman equations are just a linear system





Computing Actions from Values

Computing Actions from Q-Values

Policy Iteration

- 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)



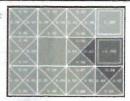
$$\pi^*(s) = \arg\max_{a} \sum_{s} T(s, a, s') [R(s, a, s') + \gamma V^*(s')]$$

- easy ac we have V* (=') for every 5.

 This is called policy extraction, since it gets the policy implied by the values
- · Given optimal values for every state, it's not immediately obvious what the optimal action is at every state (ie, optimal policy)

- Let's imagine we have the optimal q-values:
- B 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!



algo that

1) eval. policy 2) improve police

3) repeat

Problems with Value Iteration

k=0

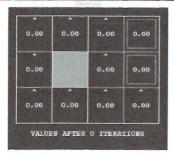
k=1

Value iteration repeats the Bellman updates:

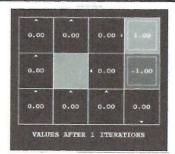
$$V_{k+1}(s) \leftarrow \max_{a} \sum_{s'} T(s, a, s') \left[R(s, a, s') + \gamma V_k(s') \right]$$

- 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

Demo value iteration (1901)



Noise = 0.2 Discount = 0.9 Living reward = 0

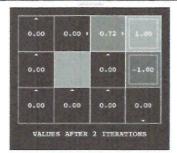


Naise = 0.2 Discount = 0.9 Living reward = 0

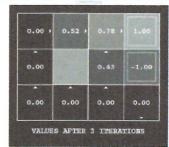
k=2

k=3

k=4

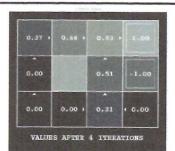


Noise = 0.2 Discount = 0.9 Living reward = 0



Noise = 0.2 Discount = 0.9 Living reward = 0

· Takes time for states for away from +1.00 + -1.00 to change arrows.



Noise = 0.2 Discount = 0.9 Living reward = 0

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
 - · Typically when # of actions is large but maximizing action doesn't change often over value iteration rounds

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 Beliman updates
 - They all use one-step lookahead expectimax fragments
 - They differ only in whether we plug in a fixed policy or max over actions

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:

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

- · Eval. is where time sowing is, and we repeat eval many to

Double Bandits

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

like this holds

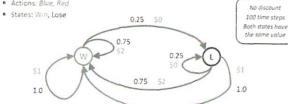
Double-Bandit MDP



given a state + action

Let's Play!

· Actions: Blue, Red

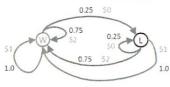


same for a particular i.e., a win or a loss is nondetermini have nondeterministic act Online Planning

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!



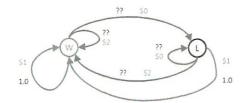


No discount 100 time steps Both states have the same value



- \$2 \$2 \$0 \$2 \$2 \$2 \$2 \$0 \$0 \$0

Rules changed! Red's win chance is different.







\$0 \$0 \$0 \$2 \$0 \$2 \$0 \$0 \$0 \$0

- That wasn't planning, it was learning!
 - Specifically, reinforcement learning
 - There was an MDP, but you couldn't solve it with just computation
 - · You needed to actually act to figure it out
- Important ideas in reinforcement learning that came up
- Exploration: you have to try unknown actions to get information
- Exploitation: eventually, you have to use what you know
- Regret: even if you learn intelligently, you make mistakes
- Sampling: because of chance, you have to try things repeatedly
 Difficulty: learning can be much harder than solving a known MDP

