

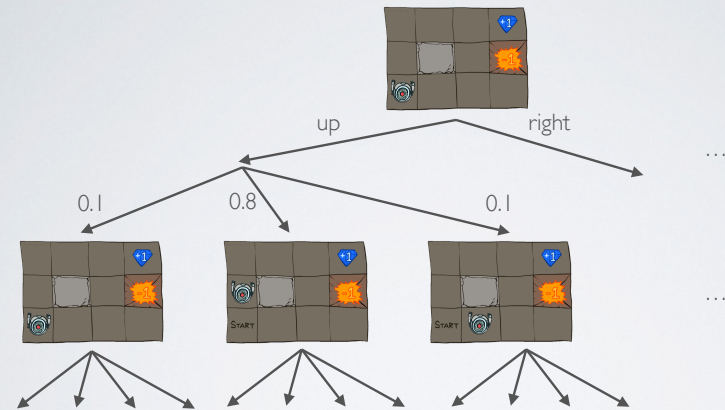
# MDP SEARCH TREES

CSE 511A: Introduction to Artificial Intelligence

Some content and images are from slides created by Dan Klein and Pieter Abbeel for CS188 Intro to AI at UC Berkeley.  
All CS188 materials are available at <http://ai.berkeley.edu>.

1

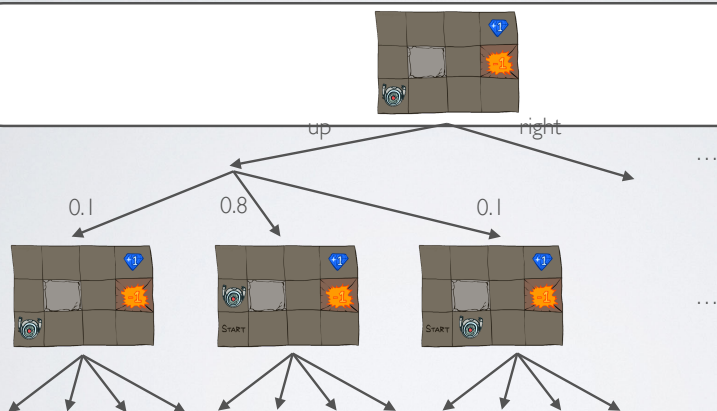
# MDP SEARCH TREE



2

# MDP SEARCH TREE

start state  $s$

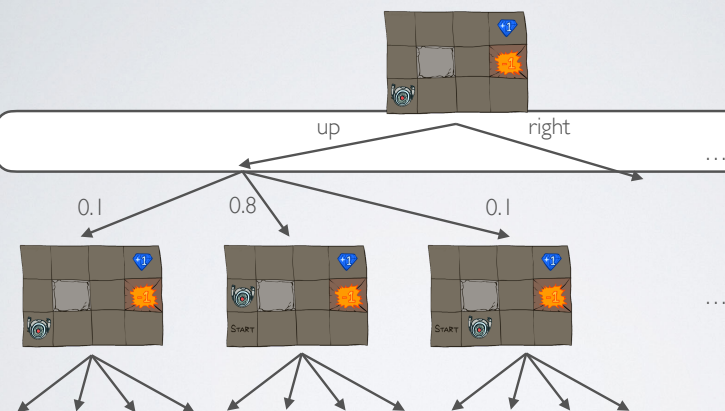


3

# MDP SEARCH TREE

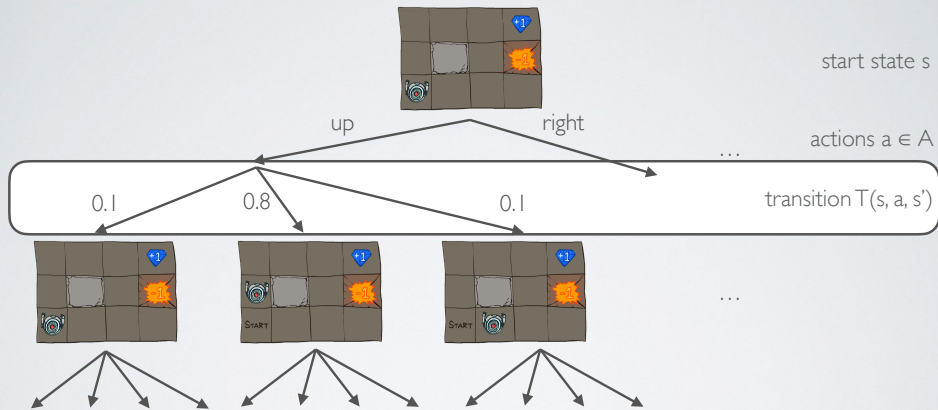
start state  $s$

actions  $a \in A$



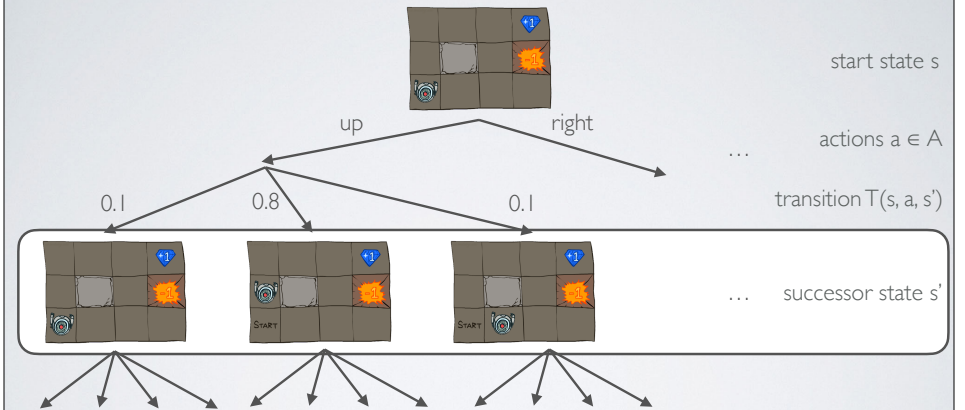
4

# MDP SEARCH TREE



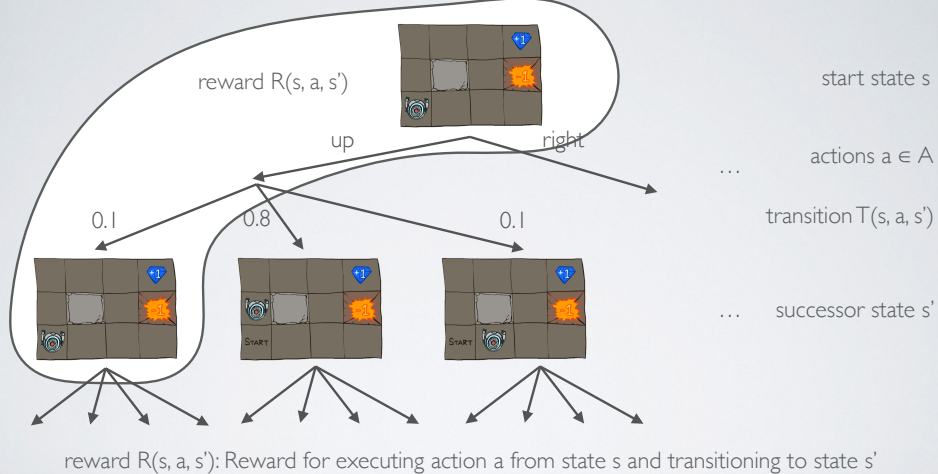
5

# MDP SEARCH TREE



6

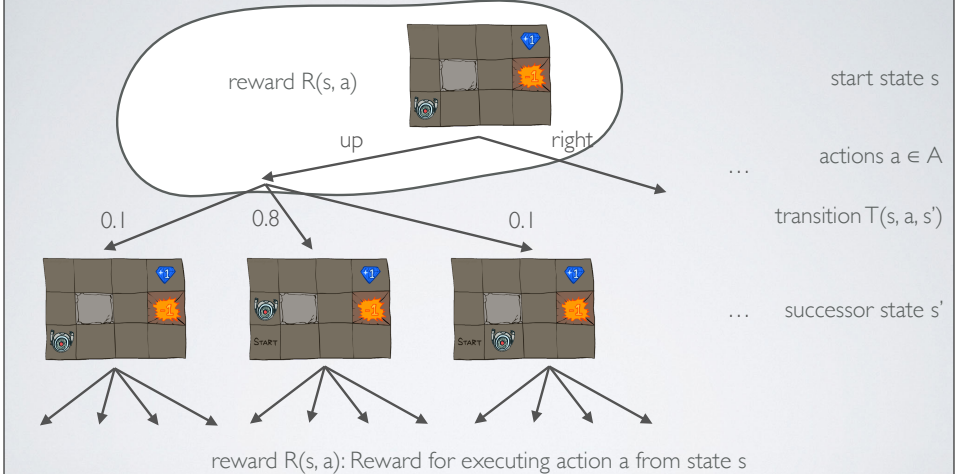
# MDP SEARCH TREE



reward  $R(s, a, s')$ : Reward for executing action  $a$  from state  $s$  and transitioning to state  $s'$

7

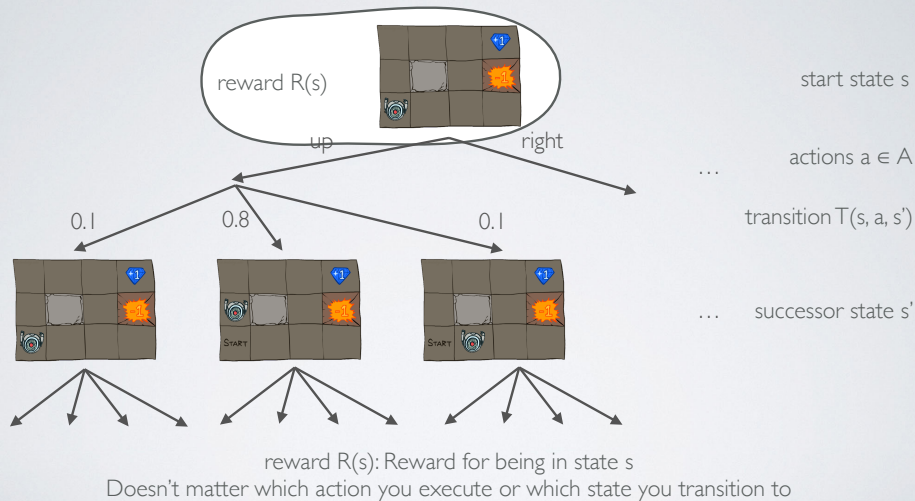
# MDP SEARCH TREE



reward  $R(s, a)$ : Reward for executing action  $a$  from state  $s$   
Doesn't matter which state you transition to

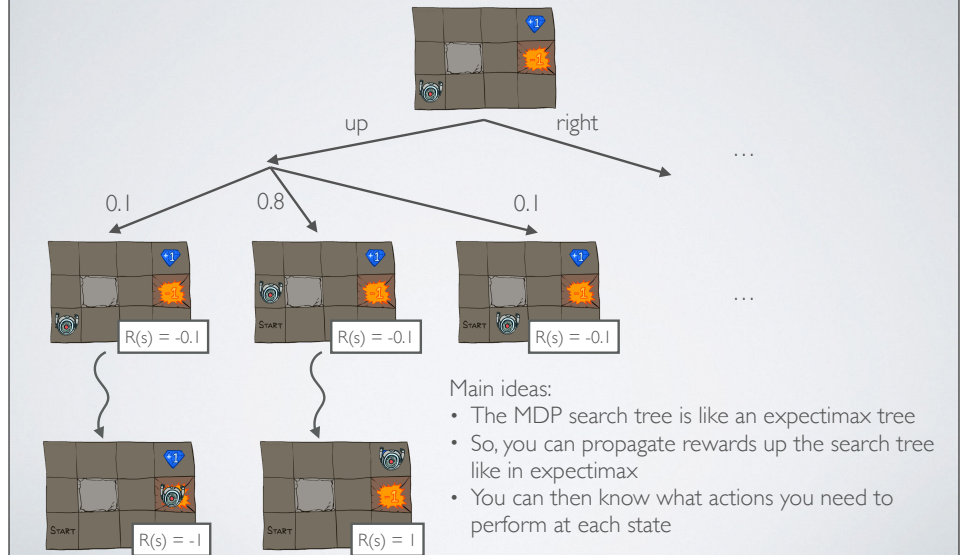
8

# MDP SEARCH TREE



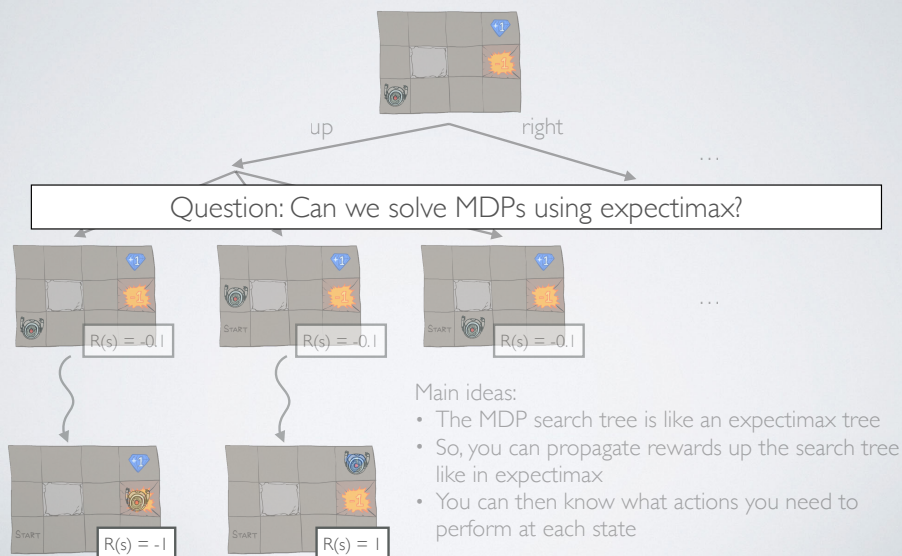
9

# SOLVING MDPs



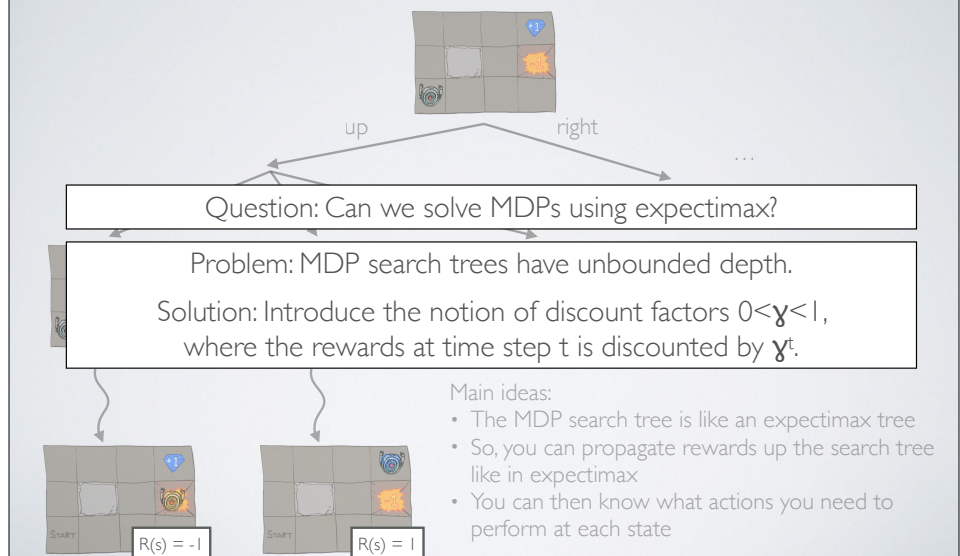
10

# SOLVING MDPs



11

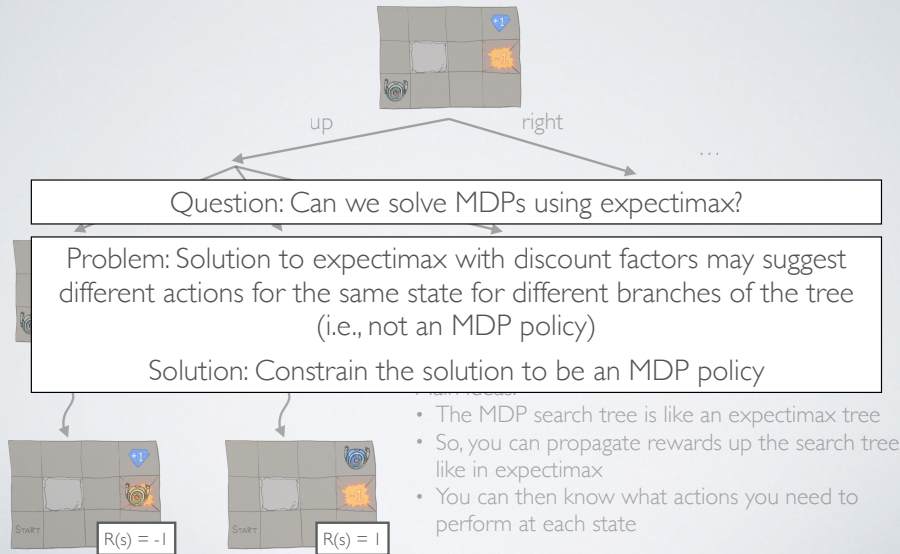
# SOLVING MDPs



12



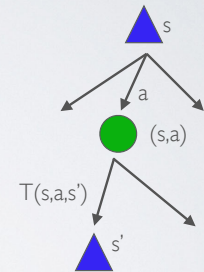
# SOLVING MDPS



13

# SOLVING MDPS

- The value (utility) of a q-state  $(s,a)$ :
  - $Q^*(s,a)$  = expected utility starting in  $s$ , taking action  $a$ , and thereafter acting optimally.
  - Important note: The action  $a$  may not be the optimal action to take at state  $s$ .
- The value (utility) of a state  $s$ :
  - $V^*(s)$  = expected utility starting in  $s$  and acting optimally
  - $V^*(s) = \max_a Q^*(s,a)$
- Optimal policy  $\pi^*$ 
  - $\pi^*(s)$  = optimal action to take at state  $s$ .
  - $\pi^*(s) = \operatorname{argmax}_a Q^*(s,a)$   
i.e., it is the action  $a$  that has the largest  $Q^*(s,a)$  over all possible actions



14

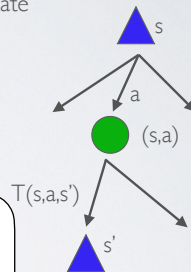
# SOLVING MDPS

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

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

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

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



15