

# CSEP 573

## Markov Decision Processes: Heuristic Search & Real-Time Dynamic Programming

Slides adapted from Andrey Kolobov and Mausam

1

## Midterm

Min

31

Mean & Median

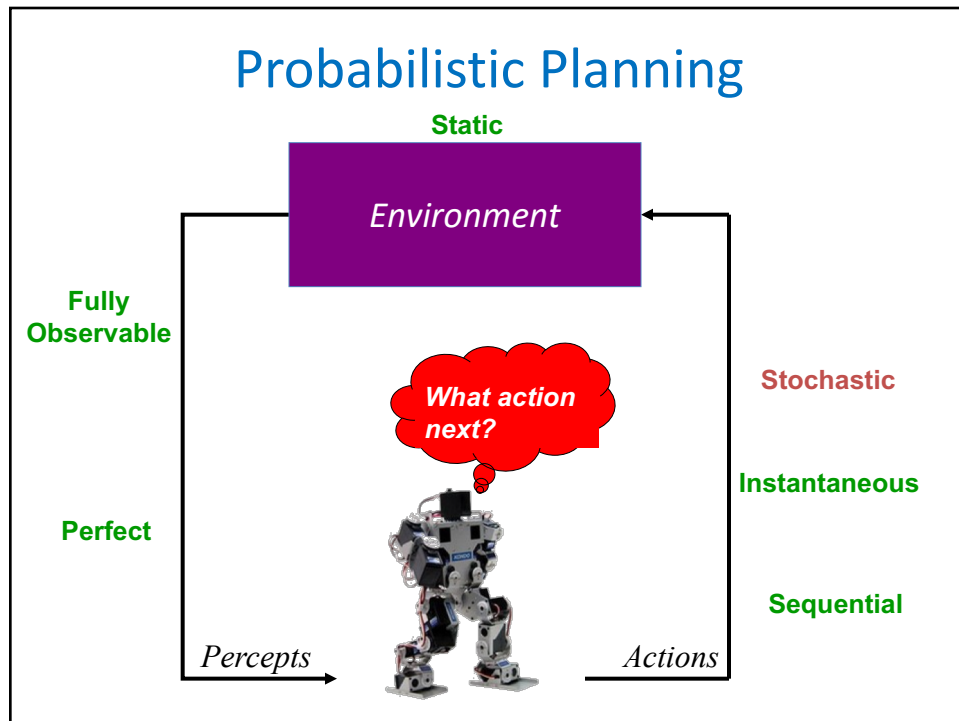
41

Max

52



Standard Dev = 5



## 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  $R(s, a, s')$

$R(s_{32}, N, s_{33}) = -0.01$   
 $R(s_{32}, N, s_{42}) = -1.01$   
 $R(s_{33}, E, s_{43}) = 0.99$

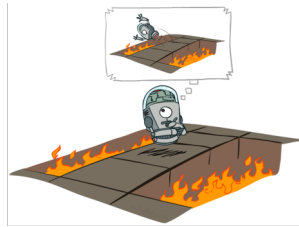
Cost of breathing

*R & T are big tables!*  
*For now, we give them to the agent*

The grid world diagram shows a 3x4 grid of cells. The bottom-left cell (row 1, column 1) is labeled "START". The bottom-right cell (row 1, column 4) contains a bomb with a red "1" and a green arrow pointing towards it. The top-right cell (row 3, column 4) contains a blue diamond with a "1". The middle-left cell (row 2, column 2) is shaded gray. The grid is indexed with rows 1, 2, 3 and columns 1, 2, 3, 4.

- We want to output the optimal policy,  $\pi^*: S \rightarrow A$ 
  - One that maximizes expected value of the discounted sequence of rewards.

## Offline vs. Online



Offline Solution  
(Planning)

Think hard; compute policy; then act



Online Learning  
(RL)

Act; learning as you go

## The Bellman Equations

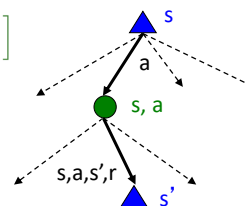
Definition of “optimal utility” via expectimax recurrence gives a simple one-step lookahead relationship amongst optimal utility values

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

These  $Q^*(s, a) = \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^*(s')]$  values are a way we can see over and over.



(1920-1984)



## Value Iteration [Bellman 57]

```

1 initialize  $V_0$  arbitrarily for each state
2  $n \leftarrow 0$ 
3 repeat
4    $n \leftarrow n + 1$ 
5   foreach  $s \in \mathcal{S}$  do
6     compute  $V_n(s)$  using Bellman backup at  $s$ 
7     compute  $\text{residual}_n(s) = |V_n(s) - V_{n-1}(s)|$ 
8   end
9 until  $\max_{s \in \mathcal{S}} \text{residual}_n(s) < \epsilon$ ;
10 return greedy policy:  $\pi^{V_n}(s) = \operatorname{argmin}_{a \in \mathcal{A}} \sum_{s' \in \mathcal{S}} T(s, a, s') [C(s, a, s') + V_n(s')]$ 

```

*iteration n*

Aka dynamic programming

Generating optimal policy assuming n time steps

In terms of optimal policy for n-1 steps

Converges to optimal policy for infinite horizon

35

## (General) Asynchronous VI

```

1 initialize  $V$  arbitrarily for each state
2 while  $\text{Res}^V > \epsilon$  do
3   select a state  $s$ 
4   compute  $V(s)$  using a Bellman backup at  $s$ 
5   update  $\text{Res}^V(s)$ 
6 end
7 return greedy policy  $\pi^V$ 

```

$$\text{Res}^V(s) = |V(s) - \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + V(s')]|$$

$$\text{Res}^V = \max_s \text{Res}^V(s)$$

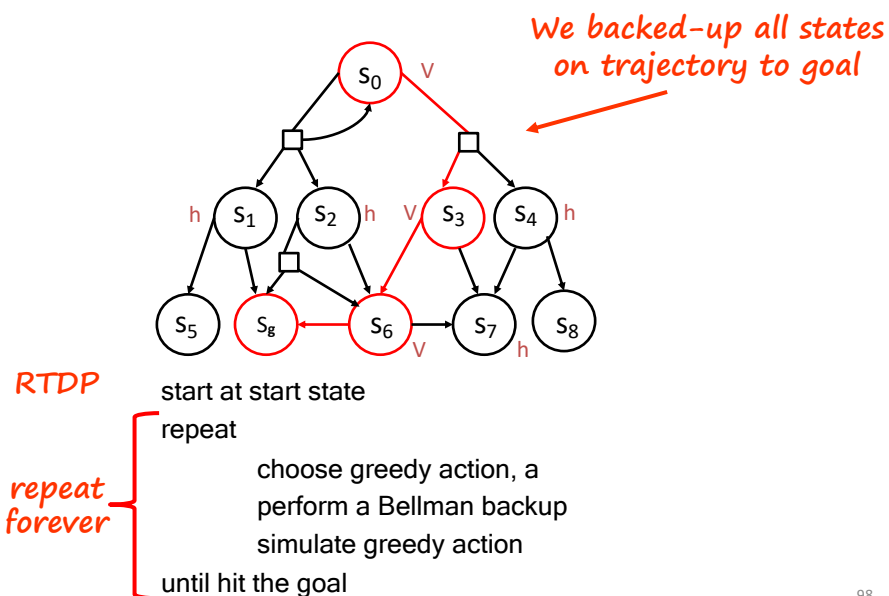
# Real Time Dynamic Programming

[Barto et al 95]

- An *instance* of asynchronous value iteration
- RTDP: repeat trials forever
  - Converges in the limit #trials  $\rightarrow \infty$
  - (doesn't test residual)
- Trial
  - simulate greedy policy starting from start state;
  - perform Bellman backup on visited states
  - stop when you hit the goal (or after N actions)
- Original Motivation
  - Agent acting in the real world (**online**)
  - But also useful as a planning algorithm (**offline**)

95

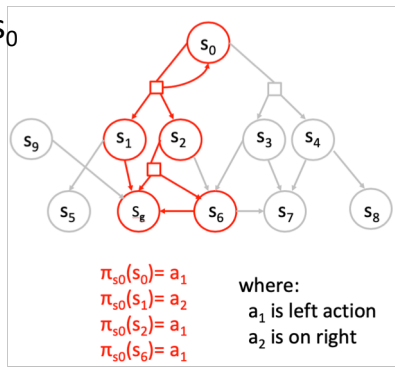
## Trial



98

## Creates a Greedy Partial Policy

- Define *greedy policy*:  $\pi^V = \operatorname{argmin}_a Q^V(s,a)$
- Define *greedy partial policy rooted at  $s_0$* 
  - Partial policy rooted at  $s_0$
  - Greedy policy  $\pi^V$
  - denoted by  $\pi^{s_0}$

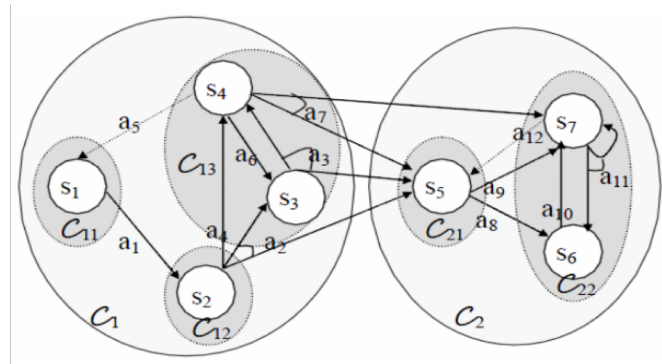


## RTDP

- Pros**
  - anytime
  - focuses attention on more probable successor states
  - can use admissible heuristic
- Cons**
  - no termination condition (see LRTDP)
  - no emphasis on highly uncertain states (see BRTDP)

# Topological Value Iteration

Another kind of asynchronous VI



129

## Backup Order Matters VI - $k=1$



Noise = 0.2  
Discount = 0.9  
Living reward = 0

$k=2$



Noise = 0.2  
Discount = 0.9  
Living reward = 0

$k=3$



Noise = 0.2  
Discount = 0.9  
Living reward = 0



## TVI: When may it perform badly?

- Highly interconnected MDP
- TVI pays overhead of computing connected components
- Reaps no benefit

133

## TVI: Convincing Experiments?

- Authors cherry-picked examples with many CCs
- And only 2 examples
- And one is artificial
- What should they have done?

134