

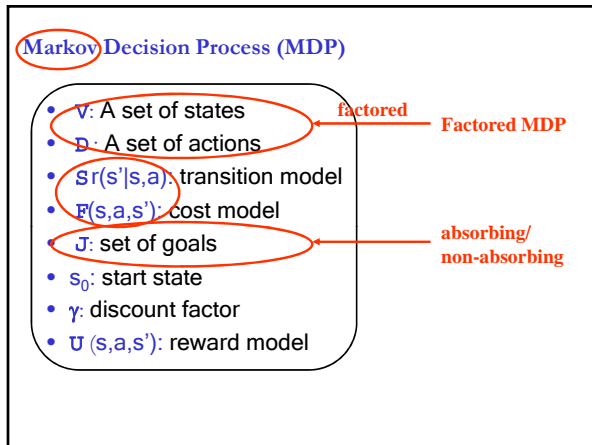
# Markov Decision Processes

## Chapter 17

Mausam

### MDP vs. Decision Theory

- Decision theory - episodic
- MDP -- sequential



### Objective of an MDP

- Find a policy  $\pi: V \rightarrow D$
- which optimizes
  - minimizes  $\left[ \begin{matrix} \text{discounted} \\ \text{or} \end{matrix} \right]$  expected cost to reach a goal
  - maximizes  $\left[ \begin{matrix} \text{discounted} \\ \text{or} \\ \text{undiscount} \end{matrix} \right]$  expected reward
  - maximizes  $\left[ \begin{matrix} \text{discounted} \\ \text{or} \\ \text{undiscount} \end{matrix} \right]$  expected (reward-cost)
- given a \_\_\_ horizon
  - finite
  - infinite
  - indefinite
- assuming full observability

### Role of Discount Factor ( $\gamma$ )

- Keep the total reward/total cost finite
  - useful for infinite horizon problems
- Intuition (economics):
  - Money today is worth more than money tomorrow.
- Total reward:  $r_1 + \gamma r_2 + \gamma^2 r_3 + \dots$
- Total cost:  $c_1 + \gamma c_2 + \gamma^2 c_3 + \dots$

### Examples of MDPs

- Goal-directed, Indefinite Horizon, Cost Minimization MDP
  - $\langle V, D, Sr, F, J, s_0 \rangle$
  - Most often studied in planning, graph theory communities
- Infinite Horizon, Discounted Reward Maximization MDP
  - $\langle V, D, Sr, U, \gamma \rangle$  **most popular**
  - Most often studied in machine learning, economics, operations research communities
- Oversubscription Planning: Non absorbing goals, Reward Max. MDP
  - $\langle V, D, Sr, J, U, s_0 \rangle$
  - Relatively recent model

### AND/OR Acyclic Graphs vs. MDPs

$C(a) = 5, C(b) = 10, C(c) = 1$

Expectimin works

- $V(Q/R/S/T) = 1$
- $V(P) = 6 - \text{action } a$

Expectimin doesn't work

- infinite loop
- $V(R/S/T) = 1$
- $Q(P,b) = 11$
- $Q(P,a) = \text{????}$
- suppose I decide to take a in P
- $Q(P,a) = 5 + 0.4 \cdot 1 + 0.6Q(P,a)$
- $\rightarrow = 13.5$

### Bellman Equations for MDP<sub>1</sub>

- $\langle V, D, S, r, F, J, s_0 \rangle$
- Define  $J^*(s)$  {optimal cost} as the minimum expected cost to reach a goal from this state.
- $J^*$  should satisfy the following equation:

$$J^*(s) = 0 \text{ if } s \in \mathcal{G}$$

$$J^*(s) =$$

### Bellman Equations for MDP<sub>2</sub>

- $\langle V, D, S, r, U, s_0, \gamma \rangle$
- Define  $V^*(s)$  {optimal value} as the maximum expected discounted reward from this state.
- $V^*$  should satisfy the following equation:

$$V^*(s) = \max_{a \in Ap(s)} \sum_{s' \in S} Pr(s'|s, a) [\mathcal{R}(s, a, s') + \gamma V^*(s')]$$