

# ■ MDP Cheat Sheet (Lecture 1)

## MDP Definition

An MDP is a 4-tuple (S, A, R, P):

- S: set of states
- A: set of actions
- $R(s,a)$  or  $R(s,a,s')$ : reward model
- $P(s,a,s')$ : transition probabilities

Markov property: Future depends only on (s,a), not the full past.

## Policies

- Deterministic:  $\pi(s) \in A$
- Randomized:  $\pi(s) \in \Delta(A)$
- Stationary: same rule forever
- Non-stationary: rule changes over time

## Objective Functions

Finite Horizon	$V^\pi(s) = E[ \sum_{t=0}^{H-1} \gamma^t r_{t+1} ]$
Infinite Horizon (Discounted)	$V^\pi(s) = E[ \sum_{t=0}^{\infty} \gamma^t r_{t+1} ]$
Infinite Horizon (Average)	$\rho^\pi(s) = \lim_{T \rightarrow \infty} (1/T) \sum_{t=0}^{T-1} E[r_{t+1}]$

## Value Functions

- $V^\pi(s)$ : expected return starting from s under  $\pi$
- $Q^\pi(s,a)$ : expected return starting from (s,a) then following  $\pi$

## Bellman Equations

Policy Evaluation:

$$V^\pi(s) = E[R(s,a) + \gamma \sum_{s'} P(s,a,s') V^\pi(s')]$$

Optimal Value:

$$V^*(s) = \max_a [ R(s,a) + \gamma \sum_{s'} P(s,a,s') V^*(s') ]$$

Optimal Q-value:

$$Q^*(s,a) = R(s,a) + \gamma \sum_{s'} P(s,a,s') \max_{a'} Q^*(s',a')$$

## Algorithms to Solve MDPs

- Value Iteration: iteratively apply Bellman optimality updates
- Policy Iteration: alternate evaluation and improvement
- Q-value Iteration: update Q directly
- Linear Programming: solve Bellman fixed point as LP