

# Markov Decision Process

$$\langle \mathcal{S}, \mathcal{A}, \mathcal{P}_{ss'}^a, \mathcal{R}_s^a, \gamma \rangle$$

**Discount Factor:**  $0 \leq \gamma \leq 1$

## State Transition Matrix

$$\mathcal{P}^a = [\mathcal{P}_{ss'}^a] \quad \mathcal{P}_{ss'}^a = \mathbb{P}[S_{t+1} = s' \mid S_t = s, A_t = a] \quad s, s' \in \mathcal{S}$$

## Policy

$$\pi(a \mid s) = \mathbb{P}[A_t = a \mid S_t = s]$$

## Reward Function

$$\begin{aligned} \mathcal{R}_s^a &= \mathbb{E}[R_{t+1} \mid S_t = s, A_t = a] = \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a \mathcal{R}_{ss'}^a \\ \mathcal{R}^a &= \left[ \mathcal{R}_{s_1}^a \quad \mathcal{R}_{s_2}^a \quad \cdots \quad \mathcal{R}_{s_{|\mathcal{S}|}}^a \right]^\top \quad \mathcal{R}_{s_i}^a \in \mathbb{R} \end{aligned}$$

## Return

$$G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} = R_{t+1} + \gamma G_{t+1}$$

## State-Value Function

$$\begin{aligned} v_\pi(s) &= \mathbb{E}_\pi[G_t \mid S_t = s] = \sum_{a \in \mathcal{A}} \pi(a \mid s) \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma v_\pi(s')] \\ v_\pi &= [v_\pi(s_1) \quad v_\pi(s_2) \quad \cdots \quad v_\pi(s_{|\mathcal{S}|})]^\top \quad v_\pi(s_i) \in \mathbb{R} \end{aligned}$$

## Action-Value Function

$$\begin{aligned} q_\pi(s, a) &= \mathbb{E}_\pi[G_t \mid S_t = s, A_t = a] \\ q_\pi(s, a) &= \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma v_\pi(s')] \quad v_\pi(s) = \sum_a \pi(a \mid s) q_\pi(s, a) \end{aligned}$$

## Reduction: MDP $\rightarrow$ Markov Reward Process (MRP)

Applying a *fixed policy*  $\pi$  defines an induced *Markov Reward Process* (MRP):  $\langle \mathcal{S}, \mathcal{P}_{ss'}^\pi, \mathcal{R}_s^\pi, \gamma \rangle$

$$\begin{aligned} \mathcal{P}_{ss'}^\pi &= \sum_{a \in \mathcal{A}} \pi(a \mid s) \mathcal{P}_{ss'}^a \quad \mathcal{R}_s^\pi = \sum_{a \in \mathcal{A}} \pi(a \mid s) \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a \mathcal{R}_{ss'}^a \\ \mathcal{P}^\pi &= [\mathcal{P}_{ss'}^\pi] \quad \mathcal{R}^\pi = \left[ \mathcal{R}_{s_1}^\pi \quad \mathcal{R}_{s_2}^\pi \quad \cdots \quad \mathcal{R}_{s_{|\mathcal{S}|}}^\pi \right]^\top \quad \mathcal{R}_{s_i}^\pi \in \mathbb{R} \end{aligned}$$

## Bellman Equation

$$v_\pi = \mathcal{R}^\pi + \gamma \mathcal{P}^\pi v_\pi \implies v_\pi = (I - \gamma \mathcal{P}^\pi)^{-1} \mathcal{R}^\pi$$

# Dynamic Programming

**Norms:** For a vector  $v \in \mathbb{R}^{|S|}$ , the  $L_p$  norm is

$$\|v\|_p = \left( \sum_i |v_i|^p \right)^{1/p}, \quad \|v\|_\infty = \max_i |v_i|$$

**Contraction:** An operator  $T$  is a  $\gamma$ -contraction if

$$\|T(u) - T(v)\|_\infty \leq \gamma \|u - v\|_\infty, \quad 0 \leq \gamma < 1$$

**Fixed Point:**  $v$  is a fixed point of  $T$  if

$$T(v) = v$$

## Bellman Expectation Operator

$$T^\pi(v) = \mathcal{R}^\pi + \gamma \mathcal{P}^\pi v \quad \|T^\pi(u) - T^\pi(v)\|_\infty \leq \gamma \|u - v\|_\infty$$

## Optimal Value Functions

$$\begin{aligned} v_*(s) &= \max_\pi v_\pi(s) = \max_{a \in \mathcal{A}} q_*(s, a) & v_* &= \max_{a \in \mathcal{A}} [\mathcal{R}^a + \gamma \mathcal{P}^a v_*] \quad (\text{element-wise}) \\ q_*(s, a) &= \max_\pi q_\pi(s, a) = \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma \max_{a'} q_*(s', a')] & &= \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma v_*(s')] \end{aligned}$$

## Optimal Policy

$$\pi_*(a \mid s) = \begin{cases} 1, & a \in \arg \max a' \in \mathcal{A} q_*(s, a') \\ 0, & \text{otherwise} \end{cases}$$

## Bellman Optimality Operator

$$T^*(v) = \max_{a \in \mathcal{A}} [\mathcal{R}^a + \gamma \mathcal{P}^a v] \quad (\text{element-wise}) \quad \|T^*(u) - T^*(v)\|_\infty \leq \gamma \|u - v\|_\infty$$

## Contraction Mapping Theorem

A  $\gamma$ -contraction  $T$  on a complete metric space has a *unique fixed point*  $v_e$ , and the iterates  $v_{n+1} = T(v_n)$  converge to  $v_e$  at a linear rate  $\gamma$ :  $\|v_{n+1} - v_e\| \leq \gamma \|v_n - v_e\|$ .

$$T^\pi(v_\pi) = v_\pi \quad (\text{iterative policy evaluation})$$

$$T^*(v_*) = v_* \quad (\text{value iteration})$$

## Policy Iteration

$$v_{\pi_k} = T^{\pi_k}(v_{\pi_k}) \quad (\text{evaluate current policy})$$

$$\pi_{k+1}(s) = \arg \max_{a \in \mathcal{A}} \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma v_{\pi_k}(s')] \quad (\text{improve policy})$$

$$\lim_{k \rightarrow \infty} \pi_k = \pi_* \quad \lim_{k \rightarrow \infty} v_{\pi_k} = v_*$$

# Model-Free Prediction

## Monte Carlo Estimation

Suppose a random variable  $X$  with unknown expected value  $\mu = \mathbb{E}[X]$ , the Monte Carlo estimate using  $N$  independent samples  $X_1, \dots, X_N$  is:

$$\hat{\mu}_N = \frac{1}{N} \sum_{i=1}^N X_i$$

$$\mathbb{E}[\hat{\mu}_N] = \mu \quad \hat{\mu}_N \rightarrow \mu \text{ as } N \rightarrow \infty$$

## Monte Carlo Prediction in RL

Estimate  $v_\pi(s)$  using returns  $G_t$  over many episodes (samples):

$$G_t = \sum_{k=0}^{T-t-1} \gamma^k R_{t+k+1} \quad (\text{where } T \text{ is the terminal time})$$

$$v_\pi(s) \approx \frac{1}{N(s)} \sum_{i=1}^{N(s)} G_t^{(i)} \quad (N(s): \text{visit count for state } s)$$

Incremental form:

$$V(S_t) \leftarrow V(S_t) + \frac{1}{N(S_t)} (G_t - V(S_t))$$

$$V(S_t) \leftarrow V(S_t) + \alpha (\text{Target} - V(S_t))$$

## Temporal Difference Learning

Use bootstrapped target for one next step instead of full episode:

$$V(S_t) \leftarrow V(S_t) + \underbrace{\alpha(R_{t+1} + \gamma V(S_{t+1}) - V(S_t))}_{\text{TD Target}}$$

The term  $\delta_t = R_{t+1} + \gamma V(S_{t+1}) - V(S_t)$  is the **TD error**.

## $n$ -Step TD and TD( $\lambda$ )

$n$ -step target:

$$G_t^{(n)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n V(S_{t+n}) \quad (G_t^{(\infty)} = G_t)$$

TD( $\lambda$ ) combines all  $n$ -step targets:

$$G_t^\lambda = (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} G_t^{(n)} \quad V(S_t) \leftarrow V(S_t) + \alpha(G_t^\lambda - V(S_t))$$

$\lambda = 0$ : TD(0) 1-step update       $\lambda = 1$ : TD(1) Monte Carlo update