# Finite Markov Decision Process
## (MDPs)

$\Big\{$ evaluative feedback

$\quad$ associative aspect ( choose different actions
$\qquad\qquad\qquad\qquad$ in different situations )

MDPs : sequential decision making

$\quad$ bandit : $\qquad q_*(a)$ w.r.t. action $a$

$\quad$ MDPs : $\quad \begin{array}{l} q_*(s,a) \text{ w.r.t. } \begin{cases} \text{action } a \\ \text{state } s \end{cases} \\ \\ \\ V_*(s) \text{ w.r.t. state } s \text{ given} \\ \qquad\qquad \text{optimal action selections} \end{array}$

Agent - Environment Interface

Action $A_t$

Agent

Environment

State $S_{t+1}$

reward $R_{t+1}$

Sequence / trajectory :

$$S_0, A_0, R_1, S_1, A_1, R_2, S_2, A_2, R_3, \ldots$$

$$p: S \times R \times S \times A \rightarrow [0, 1]$$

$$p(s', r \mid s, a) \doteq Pr\{S_t = s', R_t = r \mid S_{t-1} = s, A_{t-1} = a\}$$

$$\sum_{s' \in S} \sum_{r \in R} p(s', r \mid s, a) = 1$$

for all $s \in S$, $a \in A(s)$

"Markov Property"

"State-transition"

$$p : S \times S \times A \rightarrow [0, 1]$$

$$p(s' | s, a) \doteq Pr\{S_t = s' | S_{t-1} = s, A_{t-1} = a\}$$

$$= \sum_{r \in R} p(s', r | s, a)$$

## Expected rewards for state-action pairs

$$r : S \times A \rightarrow \mathbb{R}$$

$$r(s, a) \doteq \mathbb{E}[R_t | S_{t-1} = s, A_{t-1} = a]$$

$$= \sum_{r \in R} r \sum_{s' \in S} p(s', r | s, a)$$

## Expected rewards for state-action-next-state triples

$$r : S \times A \times S \rightarrow \mathbb{R}$$

$$r(s, a, s') \doteq \mathbb{E}[R_t | S_{t-1} = s, A_{t-1} = a, S_t = s']$$

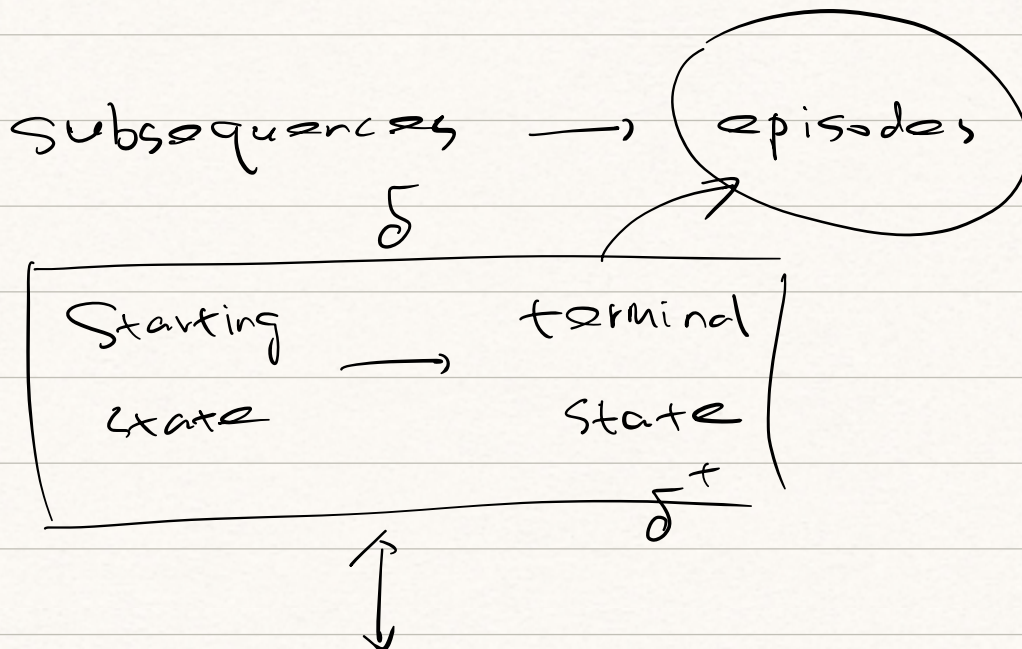$$= \sum_{r \in R} r \frac{p(s', r | s, a)}{p(s' | s, a)}$$

# Reward

$$R_t \in \mathbb{R}$$

Goal : Maximize cumulative reward in the long run

Expected return $G_t$

$$G_t \doteq R_{t+1} + R_{t+2} + \cdots + R_T$$

subsequences $\longrightarrow$ episodes

$S$

| Starting $\longrightarrow$ terminal |
| state                state |

$S^+$

$\updownarrow$

episodic task

Discounted Return :

$$G_t \doteq R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \cdots$$

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

$\gamma$ : discount rate $\qquad \gamma \in [0,1]$

$\gamma = 0$ : Agent "myopic" "immediate rewards"

$$G_t \doteq R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4} + \cdots$$

$$= R_{t+1} + \gamma (R_{t+2} + \gamma R_{t+3} + \cdots)$$

$$= R_{t+1} + \gamma G_{t+1}$$

$$G_t = \sum_{k=0}^{\infty} \gamma^k = \frac{1}{1-\gamma}$$

defined w.r.t. particular ways of acting
called policy

$$\boxed{\text{Value Functions}}$$

policy: a mapping from states to probabilities of
$\pi$

       selecting each possible action

$\pi(a|s)$ : the probability that $A_t = a$ if $S_t = s$

Value function of a state $s$ under a policy $\pi$

State-value function for policy $\pi$

$$V_\pi(s) \doteq \mathbb{E}[G_t | S_t = s]$$

$$= \mathbb{E}_\pi\left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \Big| S_t = s \right]$$

action-value function for policy $\pi$

$$q_\pi(s,a) \doteq \mathbb{E}[G_t | S_t = s, A_t = a]$$

$$= \mathbb{E}_\pi\left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \Big| S_t = s, A_t = a \right]$$

$V_\pi(s)$

$$\doteq \mathbb{E}_\pi [G_t | S_t = s]$$

$$= \mathbb{E}_\pi [R_{t+1} + \gamma G_{t+1} | S_t = s]$$

$$= \sum_a \pi(a|s) \sum_{s'} \sum_r p(s', r | s, a) [r + \gamma \mathbb{E}_\pi [G_{t+1} | S_{t+1} = s']]$$

$$= \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) [r + \gamma V_\pi(s')]$$

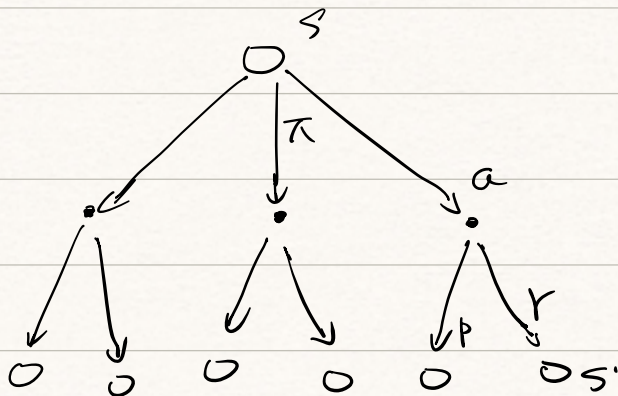Bellman equation for $V_\pi$

$$V_\pi(s) = \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) [r + \gamma V_\pi(s')]$$

transfer value information back to a state (or a
state-action pair) from its successor states

(or state-action pairs)     "Bootstrap"

Backup diagram for $V_\pi$

# Optimal Policies / Optimal Value Functions

Optimal policy: $\pi_*$

optimal state-value function $V_*$

$$V_*(s) \doteq \max_\pi V_\pi(s)$$

optimal action-value function $q_*$

$$q_*(s,a) \doteq \max_\pi q_\pi(s,a)$$

$$= \mathbb{E}\left[R_{t+1} + \gamma V_*(S_{t+1}) \mid S_t = s, A_t = a\right]$$

## Bellman Optimality Equation

$$V_*(s) = \max_{a \in A(s)} q_{\pi_*}(s,a)$$

$$= \max_a \mathbb{E}_{\pi_*}\left[G_t \mid S_t = s, A_t = a\right]$$

$$= \max_a \mathbb{E}_{\pi_*}\left[R_{t+1} + \gamma G_{t+1} \mid S_t = s, A_t = a\right]$$

$$\boxed{V_*(s) = \max_a \mathbb{E}\left[R_{t+1} + \gamma V_*(S_{t+1}) \mid S_t = s, A_t = a\right]}$$

$$V_*(s) = \max_a \sum_{s',r} P(s', r \mid s, a)\left[r + r V_*(s')\right]$$

$$q_*(s,a) = E\left[R_{t+1} + \gamma \max_{a'} q_*(S_{t+1}, a') \mid S_t = s, A_t = a\right]$$

$$= \sum_{s',r} P(s', r \mid s, a)\left[r + \gamma \max_{a'} q_*(s', a')\right]$$

Bellman optimality equation for $q_*$