# Introduction to MDP

Tuesday, February 11, 2020     11:37 PM

## Reinforcement Learning
### Formulating RL problem as Markov Decision Process (MDP)

**Markov Property**

A state $S_t$ is Markov, if and only if:
$$P[S_t|S_{t-1}] = P[S_t|S_1, S_2, \dots, S_{t-1}]$$

$$P\left(S_t \mid S_{t-1}\right) = P\left(S_t \mid S_1, S_2 \dots S_{t-1}\right)$$

**Markov Process**

$$\langle S, P \rangle$$

$$S_1, S_2 \dots S_n$$

$$r \ C_1 \qquad P(C_2|C_1) = 0.5$$
$$1 C_1 = 0.4 \qquad 0.9$$

$$\longrightarrow C_1, C_2, S$$

$\begin{bmatrix} C1 \\ C2 \\ C3 \\ Pa \\ Pu \\ S \\ R \end{bmatrix}$   $P(C_2|C_1) = 0.5$
$P(Pu|C_3) = 0.4$

$\Rightarrow C1, C2, S$

$\rightarrow C1, C2, C3, Pu, C2, C3, Pa$



**Markov Reward Process**

$\langle S, P, R \rangle$

## Markov Decision Process

$\rightarrow \langle S, P, R, A \rangle$

$\rightarrow$ Facebook
R = -1

$\rightarrow$ $S_1$, Study, $S_2$, Study,
$S_3$, Study, $S_4$.



$\rightarrow$ $S_1$, Facebook, $S_5$,
Facebook, $S_5$,
Quit, $S_1$, Study, $S_2$,
Sleep, $S_4$.

*(diagram labels)*

- $S_5$
- Facebook R = -1 (self-loop)
- Quit R = 0
- Facebook R = -1
- Sleep R = 0
- $S_4$
- $S_1$ — Study R = -2 → $S_2$ — Study R = -2 → $S_3$ — Study R = +10
- Pub R = +1
- 0.2, 0.4, 0.4
- $S_6$

## Reinforcement Learning
### Understanding the components of MDPs

**Deterministic and Stochastic Processes**

$$f(x) \to y$$

$$g(x) \to [P(y_1), P(y_2) \cdots P(y_n)]$$

$$g(x) \to \boxed{y_i}$$

## Components of MDP

$< S, A, R, P >$
$S = Finite\ set\ of\ state$
$A = Finite\ set\ of\ actions$
$R = Finite\ set\ of\ all\ rewards$
$P = Environment\ dynamics\ function$

## Understanding Environment Dynamics function

$p(s'|s,a) = \Pr\{S_t = s'|S_{t-1} = s, A_{t-1} = a\}$
$\sum_{s'} p(s'|s,a) = 1, \qquad \forall s \in S, a \in A(s)$

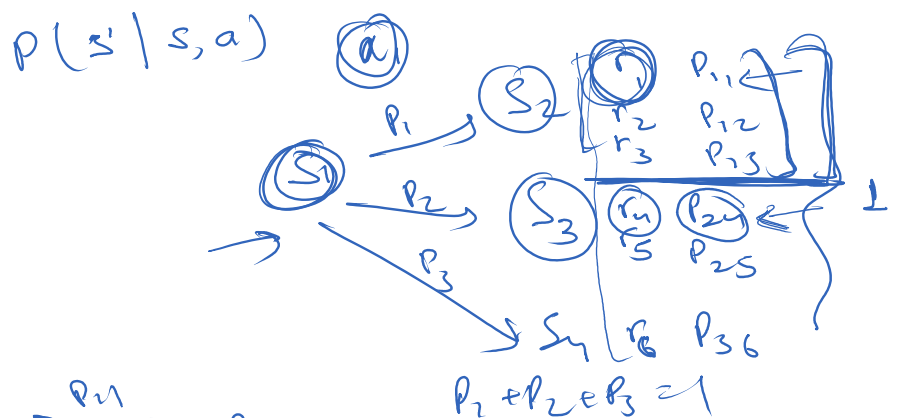$p(s',r|s,a) = \Pr\{S_t = s', R_t = r|S_{t-1} = s, A_{t-1} = a\}$
$\sum_{s'} \sum_{r} p(s',r|s,a) = 1, \qquad \forall s \in S, a \in A(s)$

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

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

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

$p\left(s' \mid s, a\right)$ (a)

$P_1 \to S_2 \to \begin{array}{cc} r_1 & P_{11} \\ r_2 & P_{12} \\ r_3 & P_{13} \end{array}$

$(S_1)$ $P_2 \to (S_3)\ (r_4)\ (P_{24}) \begin{array}{c} r_5 \\ P_{25} \end{array}$

$P_3 \to S_4\ r_6\ P_{36}$

$P_1 + P_2 + P_3 = 1$

$\frac{P_{11}}{P_{11} + P_{12} + P_{13}}$

$\Rightarrow P\left(s', r \mid s, a\right)$

$P\left(S_2 \mid s, a_1\right) = P_1$
$= P_{11} + P_{12} + P_{13}$

$r\left(s, a_1\right) = r_1 \times P_{11} + r_2 \times P_{12} + \cdots\ r_6 \times P_{36}$

$r\left(s_1, a_1, S_2\right) =$

$$\left\{ \begin{array}{l} \text{Choices (Actions)} \\ \text{States.} \\ \text{Rewards} \end{array} \right.$$

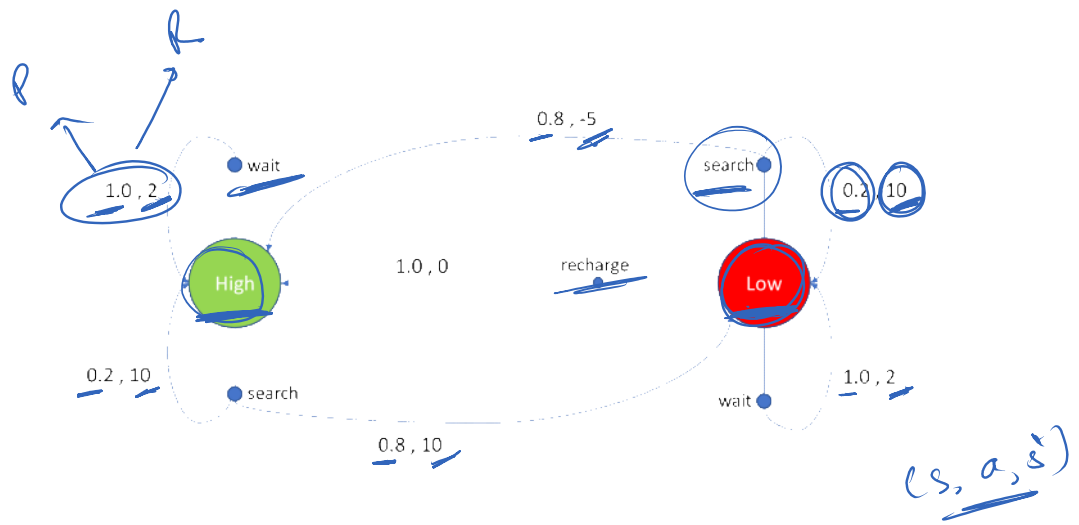# Reinforcement Learning
## MDP example

**Recycling Robot example for MDP**

$S = \{low, high\}$

$A = \{search, wait, recharge\}$

$R = \{R_{search}, R_{wait}, R_{rescued}\}$

$A(high) = \{search, wait\}$

$A(low) = \{search, wait, recharge\}$.

0.8 , -5

search

0.2 , 10

1.0 , 2

wait

1.0 , 0

recharge

Low

High

0.2 , 10

search

1.0 , 2

wait

0.8 , 10

$(s, a, s')$

| $s$ | $a$ | $s'$ | $r(s,a,s')$ | $p(s'\|s,a)$ | $p(s',r\|s,a)$ |
|---|---|---|---|---|---|
| low | search | low | 10 | 0.2 | 0.2 |
| low | search | high | -5 | 0.8 | 0.8 |
| low | wait | low | 2 | 1.0 | 1.0 |
| low | recharge | high | 0 | 1.0 | 1.0 |
| high | wait | high | 2 | 1.0 | 1.0 |
| high | search | high | 10 | 0.2 | 0.2 |
| high | search | low | 10 | 0.8 | 0.8 |

## Reinforcement Learning
### Episodic and continuing tasks



Episodic

$S_0 \xrightarrow{+1} S_1 \xrightarrow{+1} S_2 \xrightarrow{+1} S_3$

Continuing

$S_0 \xrightarrow{+1} S_1 \xrightarrow{+1} S_2 \xrightarrow{+1} S_3 \xrightarrow{0} A$

$S_0, S_1, S_2, S_3, A, A, A, \cdots$

$+1 \quad +1 \quad +1 \qquad 0 \quad 0 \quad 0 \quad 0$

Reinforcement Learning
Rewards and Returns

**Rewards and returns**

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

*With discounting,*

$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \cdots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$

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



$t \longrightarrow R_t$

$\text{return} \quad , \quad G_t = R_{t+1} + R_{t+2} + \cdots + R_T = +6$

C1  C2  C3  Passed  Sleep.

$G(G_1) = -2 + (-2) + 10 + 0$

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

$\gamma \in [0, 1]$

$\gamma < 1$

$R_t = 1$

$$G_t = 1 + \gamma + \gamma^2 + \gamma^3 + \cdots$$

$$= \frac{1}{1 - \gamma}$$

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \cdots \quad (\gamma^{k-1}) R_{t+k}$$

$$= R_{t+1} + \gamma \underbrace{\left( R_{t+2} + \gamma R_{t+3} + \gamma^2 R_{t+4} \right)}_{G_{t+1}}$$

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

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

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

# Reinforcement Learning
## Policy

$$\Pi(s)$$

Deterministic     $\Pi(s_1) \longrightarrow a$

Stochastic        $\Pi(s_1) \longrightarrow P(a_1), P(a_2), P(a_3) \ldots$

$\longrightarrow \quad \Pi(a_1|s) \longrightarrow y_1 \longrightarrow P(a_1)$

$\longrightarrow \quad \Pi(a_2|s) \longrightarrow y_2 \longrightarrow P(a_2)$



Actions

25%

$25\% \Leftarrow \downarrow \rightarrow 25\%$

$\longrightarrow \quad 25\%$

$$\kappa \left( \uparrow | A \right) = 95\%$$

$$\kappa \ 100\%$$

| A | U |
|---|---|

$$\pi \left( \uparrow | A \right) = 25\%$$

$$\pi \left( \leftarrow | A \right) = 25\%$$

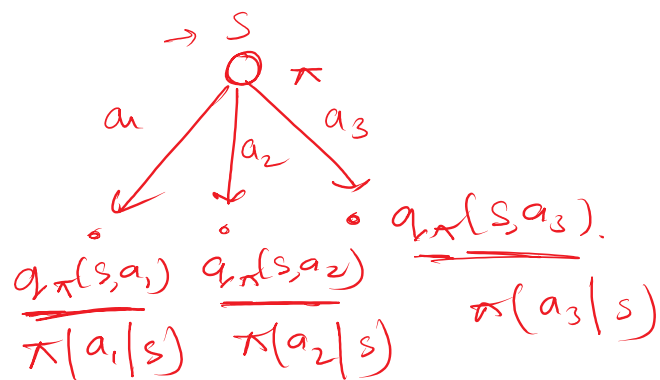# Reinforcement Learning
## State and Action Value functions

**State Value function**

$v_\pi(s)$ .

$$v_\pi(s) = \mathbb{E}_\pi[G_t | S_t = s] = \mathbb{E}_\pi[\sum_{k=0}^{T} \gamma^k R_{t+k+1} | S_t = s], \quad \forall \ s \in S$$

State value

$$G_t \to R \to A \to \pi$$

**Action Value function**

$q_\pi(a, s)$

$$q_\pi(s, a) = \mathbb{E}_\pi[G_t | S_t = s, A_t = a] = \mathbb{E}_\pi[\sum_{k=0}^{T} \gamma^k R_{t+k+1} | S_t = s, A_t = a], \quad \forall \ s \in S, a \in A(s)$$

$$\to \underset{Q}{S} \pi$$

$v_\pi$ using $q_\pi$

$v_\pi(s)$



$$v_\pi(s) = \pi(a_1|s) \cdot q_\pi(s,a_1) + $$
$$\pi(a_2|s) \cdot q_\pi(s,a_2) + $$
$$\pi(a_3|s) \cdot q_\pi(s,a_3).$$

$$v_\pi(s) = \sum_a \pi(a|s) \, q_\pi(s,a).$$

$q_\pi$ from $v_\pi$



$$q_\pi(s,a) = p(s_1,r_1|s,a)[r + \gamma \, v_\pi(s_1)] + $$
$$p(s_2,r_2|s,a)[r + \gamma \, v_\pi(s_2)] + $$
$$p(s_3,r_3|s,a)[r + \gamma \, v_\pi(s_3)]$$

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

# Reinforcement Learning
## Bellman's Equations

**Bellman's Equation for state values**

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

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

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

$$= \mathbb{E}_\pi[R_{t+1} | S_t = s] + \gamma \mathbb{E}_\pi[\mathbb{E}_\pi[G_{t+1} | S_{t+1} = s'] | S_t = s]$$

$$= \mathbb{E}_\pi[R_{t+1} | S_t = s] + \gamma \mathbb{E}_\pi[v_\pi(S_{t+1} = s') | S_t = s]$$

$$= \mathbb{E}_\pi[R_{t+1} + \gamma v_\pi(S_{t+1} = s') | S_t = s]$$

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

$$E[X] = E\left[E[X|Y]\right]$$
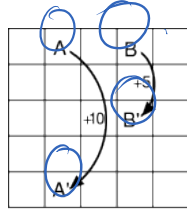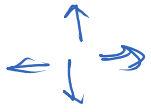
$$v_\pi(s') = E_\pi[G_{t+1} | S_{t+1} = s']$$

$$\pi(a_3|s) \cdot p(r_3, s_3 | s, a)$$

$$[r_3 + \gamma v_\pi(s')]$$

$\pi(\cdot|s) = \frac{1}{4}$.

$\gamma = 0.9$

**Figure (a):** gridworld with A, B, A', B' cells, +10, +5, 25%

**Figure (b):**

| 3.3 | 8.8 | 4.4 | 5.3 | 1.5 |
| 1.5 | 3.0 | 2.3 | 1.9 | 0.5 |
| 0.1 | 0.7 | 0.7 | 0.4 | -0.4 |
| -1.0 | -0.4 | -0.4 | -0.6 | -1.2 |
| -1.9 | -1.3 | -1.2 | -1.4 | -2.0 |

Actions

$$\frac{1}{4} \times 1 \left[0 + 0.9(2.3)\right] + \frac{1}{4} \times 1\left[0 + 0.9(0.7)\right] + \frac{1}{4}\left[0 + 0.9 \times (0.4)\right] + \frac{1}{4}\left[0 + 0.9 \times (0.4)\right]$$

$$= \quad 0.67 \sim 0.7$$

**Bellman's Equation for action values**

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

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

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

$$= \mathbb{E}_\pi[R_{t+1}|S_t = s, A_t = a] + \gamma\mathbb{E}_\pi[\mathbb{E}_\pi[G_{t+1}|S_{t+1} = s']|S_t = s, A_t = a]$$
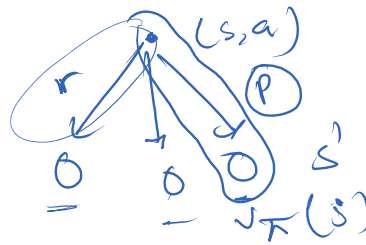
$$= \mathbb{E}_\pi[R_{t+1}|S_t = s, A_t = a] + \gamma\mathbb{E}_\pi[v_\pi(S_{t+1} = s')|S_t = s, A_t = a]$$

$$= \mathbb{E}_\pi[R_{t+1} + \gamma v_\pi(S_{t+1} = s')|S_t = s, A_t = a]$$
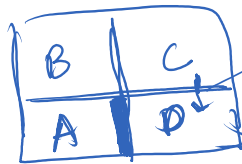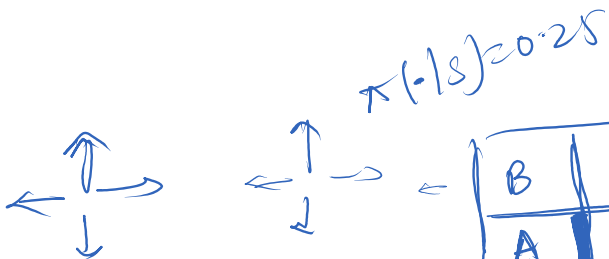
$$= \sum_{s',r} p(s', r|s, a)[r + \gamma v_\pi(s')]$$

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

$$v_\pi(s) = \sum_a \pi(a|s) q_\pi(s,a)$$

$(s,a)$

$p(s', r|s, a)$

$[r + \gamma v_\pi(s')]$

$\pi(\cdot|s) = 0.25$

$+5 \qquad \gamma = 0.7$



$$V(A) = \frac{1}{4} \times 1\left[0 + 0.7 \cdot V(B)\right] + \frac{3}{4} \times 1\left[0 + 0.7\, V(A)\right]$$

$$V(B) = \frac{1}{4} \times 1\left[0 + 0.7 \cdot V(A)\right] + \frac{1}{4} \times 1\left[0 + 0.7\, V(C)\right]$$

$$+ \quad \cdots \left[0 + 0.7\, V(B)\right]$$

$$V(B) = \frac{1}{4} \times 1 \left[ 0 + 0.7 \cdot V(A) + \frac{2}{4} \quad \wedge \quad \smile \right.$$
$$+ \frac{2}{4} \times 1 \left[ 0 + 0.7 \; V(B) \right]$$

$$V(C) = \frac{1}{4} \times 1 \left[ 0 + 0.7 \times V(B) \right] + \frac{1}{4} \times 1 \left[ 5 + 0.7 \; V(D) \right]$$
$$+ \frac{2}{4} \times 1 \left[ 0 + 0.7 \times V(C) \right]$$

$$V(D) = 0$$

# Reinforcement Learning
## Optimality

**Optimal Policy**

$\pi_1 \quad \pi_2$

$\pi_1 \geq \pi_2$ iff $v_{\pi_1}(s) \geq v_{\pi_2}(s) \; \forall s \in S.$



State values / States

$\pi_2 \ngeq \pi_1$

$\pi_1 \ngeq \pi_2$

$\pi_3 \geq \pi_1$

$\pi_3 \geq \pi_2$

$\pi_4 \geq \pi_1, \pi_2, \pi_3$

$$v_*= \max_\pi v_\pi(s) \; \forall s \qquad \pi_*$$

$$q_*(s,a) = \max_\pi q_\pi(s,a), \; \forall s,a$$

## Bellman optimality equations

$$v_*(s) = \max_a \sum_{s',r} p(s',r|s,a)[r + \gamma v_*(s')]$$

$$q_*(s,a) = \sum_{s',r} p(s',r|s,a)[r + \gamma \max_a q_*(s',a')]$$



$$v_*(s) = \max_a q_{\pi_*}(s,a)$$



a) gridworld

b) V*

c) π*