# Solve MRP

# MRP

① S ② $P(s, s')$, ③ $R$ ④ $r$.

⑤ $T = \inf \{t \geq 0, S_t \in \partial S\}$

## Episode

$$\omega = \{S_0, R_1, S_1, \text{----} \quad R_T, S_T, \colorbox{yellow}{$R_{T+1}$}\}$$

Terminal Reward

## Gain

$$G_t = R_{t+1} + r R_{t+2} + \cdots + r^{T-t} R_{T+1}$$

## Value

$$V(s) = \mathbb{E}\left[G_t \mid S_t = s\right]$$

# Bellman

$$\begin{cases} v(s) = \mathbb{E}[R_{t+1} + \gamma v(S_{t+1}) \mid S_t = s] & \forall s \notin \partial S \\ v(s) = R(s) & \forall s \in \partial S \end{cases}$$

## To solve Bellman

$$\begin{cases} \textcircled{1} \text{ linear Algebra} \\ \textcircled{2} \text{ PP} \\ \textcircled{3} \text{ MC} \\ \textcircled{4} \text{ TD} \end{cases}$$

$$v(s) = R(s) + \gamma \sum_{s'} P(s, s') v(s')$$

$$\begin{cases} v(s) = \mathbb{E}[R_{t+1} + \gamma v(S_{t+1}) \mid S_t = s]. & \forall s \notin \partial S \\ \\ v(s) = R(s) & \forall s \in \partial S \end{cases}$$

① Lin Alg.

    ⊛ Denote $\{v(s) : s \in S\}$ be a column vector
                 $\{R(s) : s \in S\}$ - - - - - - -

    ⊛ write Bellman

$$V = R + \gamma P V$$

$$V = (I - \gamma P)^{-1} R$$

② DP

$$V^{(n+1)} \longleftarrow R + \gamma P V^{(n)}$$

③ <span style="color:red">**MC1 (first visit)**</span>

To compute $V(s)$, use

$$V(s) = \mathbb{E}[G_0 \mid S_0 = s]$$

$$\approx \frac{1}{n} \sum_{i=1}^{n} G_0^{(i)}$$

<span style="color:red">**Algo**</span>

Tot = 0

for $i = 1 \cdots n$ do

    Generate $W_i = \{S_0 \overset{=s}{,} R_1, S_1, \cdots R_T, S_T, R_{T+1}\}$

    Compute $G \leftarrow R_1 + \gamma R_2 + \cdots + \gamma^T R_{T+1}$

    Tot $\leftarrow$ Tot + G

return Tot$/n$

**(B'')**  MC 2 — (every visit)

Algo (To compute $V(s)$)

$Tot = 0$

For $i = 1 \cdots n$ do

$$w^{(i)} = \{s_0, R_1, s_1 \cdots R_T, S_T, R_{T+1}\}$$

$$t^i(0) = \inf\{t \geqslant 0, S_t = s\}$$

$$t^i(j+1) = \inf\{t > t^i(j), S_t = s\} \wedge T$$

$$M^i = \#\{0 \leq t < T : S_t = s\}$$

$$V(s) = \frac{1}{\sum\limits_{i=1}^{n} M_i} \sum\limits_{i=1}^{n} \sum\limits_{j=0}^{M^i} G(t^i(j))$$

Let $(x_1, x_2 \cdots)$ be a seq.

$$\mu_k = \frac{1}{k} \sum_{i=1}^{k} x_i$$

Then $\mu_k = \mu_{k-1} + \frac{1}{k}(x_k - \mu_{k-1})$

Pf (skip)

Algo    $\mu = 0$.

for $k = 1, 2 \cdots n$ do    (learning rate)

Generate    $x_k$    $\to \alpha$

$\mu \leftarrow \mu + \frac{1}{k}(x_k - \mu)$

Return $\mu$.

**MC**

**algo** : To compute $\{V(s) : s \in S\}$ for MRP

Init $V(s) = 0 \quad \forall s \in S$

for $i = 1 \cdots n$ do

    Generate $\cdot \omega = \{S_0, R_1, S_1, \cdots R_T, S_T, R_{T+1}\}$

    $G \leftarrow R_{T+1}$

    for $t = T-1 \cdots 0$ do

        $G \leftarrow R_{t+1} + \gamma G$ .

        $V(S_t) \leftarrow V(S_t) + \alpha \left( G - V(S_t) \right)$

Return $V$.

"MC2"

$$v(S_t) \leftarrow v(S_t) + \alpha \left( G_t - v(S_t) \right)$$

"TD(0)"

$$v(S_t) \leftarrow v(S_t) + \alpha \left( \underbrace{R_{t+1} + \gamma v(S_{t+1})}_{\text{TD-target}} - v(S_t) \right)$$

$\delta_t$ : TD error

<u>pros</u>    No need for a complete episode
for update the value

**Algo** Compute V by TD(0)

Init $V \equiv 0$

for $i = 1 \cdots n$ do
    Generate $S$,
      while $S \notin \partial S$ do
        Generate $S' \sim P(S, S')$

$$V(s) \leftarrow V(s) + \alpha \left( R(s) + \gamma V(s') - V(s) \right)$$

$$S \leftarrow S'$$

Return V