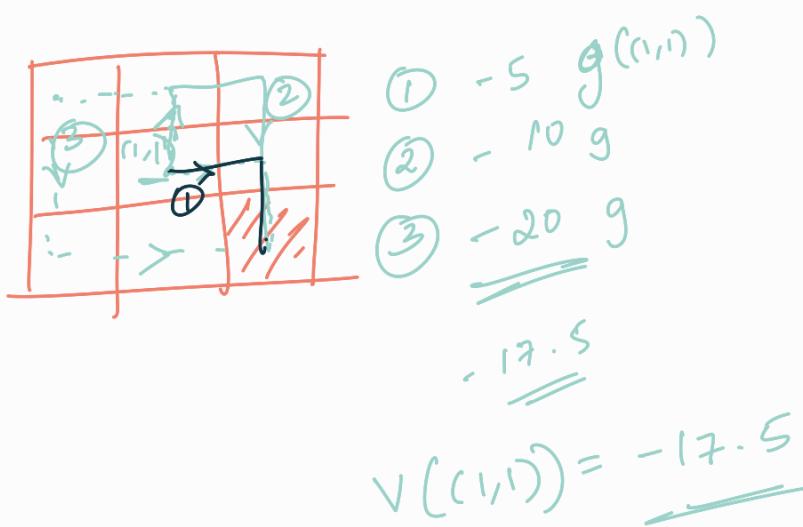


Monte Carlo

- Agent learns from experience
- Complete knowledge of env is not needed
- Agent can learn from a model of the env., but does not need to know the probability distributions of all transitions. $P(s'|s, a)$

Ex.



$$V(s) =$$

3.2	13.2	-15.1
10.5	-17.5	-1.7
-4.5	-5.6	10.

- ① policy evaluation (Monte Carlo prediction)
- ② policy improvement (Monte Carlo control)

Monte Carlo Prediction

$$V(S) = \frac{1}{C(S)} \sum_{m=1}^M \sum_{t=0}^{T-1} \mathbb{I}(S_t^m = S) g_t^m$$

Example:-

T	1	2	3	4	5	6	Given, $\pi(a s)$
S	s_0	s_1	s_0	s_0	s_1	s_0	$\pi(a s)$
A	a_0	a_0	a_1	a_0	a_1	-	
R	-	s	10	15	-3	2	

(1 episode)

$$V_{\pi}(s_0) = ?$$

$$\begin{aligned}
 ① \rightarrow G_t &= R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \\
 &= \underbrace{R_{t+1} + \gamma G_{t+1}}
 \end{aligned}$$

$$② \rightarrow V_{\pi}(s) = \underbrace{E[G_t | S_t = s]}$$

$$\gamma = 1,$$

Sol:-

$$g_6 = 0$$

$$g_5 = R_6 + \gamma g_6 = 2 + 1(0) = 2$$

$$g_4 = R_5 + \gamma g_5 = -3 + 2 = -1$$

$$g_3 = 15 - 1 = 14$$

$$g_2 = 10 + 14 = 24$$

$$g_1 = 5 + 24 = 29$$

$$\begin{aligned}
 \underline{V_{\pi}(s_0)} &= \frac{g_1 + g_3 + g_4 + g_6}{4} \\
 &= \frac{29 + 14 + (-1) + 0}{4} \\
 &= \frac{42}{4} = 10.3
 \end{aligned}$$

Iterative Calculation of mean

$$\rightarrow \frac{x_0 + x_1 + x_2 + \dots + x_n}{n} = \bar{x}_n$$

$$\begin{aligned}
 \bar{x}_n &= \bar{x}_{n-1} + \frac{1}{n} (x_n - \bar{x}_{n-1}) \\
 &= \frac{n \bar{x}_{n-1} + x_n - \bar{x}_{n-1}}{n}
 \end{aligned}$$

$$= \frac{\bar{x}_{n-1}(n-1) + x_n}{n}$$

$$= \frac{x_0 + x_1 + x_2 + \dots + x_{n-1} \times (n-1) + x_n}{n}$$

$$= \bar{x}_n$$

Monte Carlo prediction (policy evaluation)

→ Given $\pi(a|s)$, $\underline{V(s)} \approx \underline{V_\pi(s)}$

Initialise

$$V(s) = \text{random}$$

loop many times, (ideally loop forever)
so choose the start state randomly.
Simulate one episode under $\pi(a|s)$
→ $(s_0, a_0, r_0), (s_1, a_1, r_1), s_2, \dots \dots \dots$
= history

② Update $\underline{V(s)}$ based on the episode.

$$G_t = 0$$

for $t = T-1, T-2, \dots, 0$

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

$$\rightarrow \underline{V(s_t)} = \underline{V(s_t)} + \frac{1}{N(s_t)} \left(G_t - \underline{V(s_t)} \right)$$

state
 s_t is visited

→ For 4×4 Gridworld, evaluate
a random policy (actions are chosen
randomly). → code the MC prediction
algorithm.

$$* V_*(s) = \max_{a \in A} \left[\sum_{s' \in S} P(s'|s, a) \left[\underbrace{r(s, a, s')}_{\text{Bellman optimality equation for state}} + \gamma V_*(s') \right] \right]$$

Bellman's optimality equation for state
value function

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

Monte Carlo Control

we maximise the state action value function.

$$S = \{s_0, s_1, s_2\}$$

$$A = \{a_0, a_1\}$$

$q_{\pi}(s, a) \rightarrow Q\text{-table}$			
	a_0	a_1	
s_0	1.25	1.75	a_1
s_1	1.3	1.4	a_1
s_2	1.3	1.0	a_0

ϵ -greedy

$$0 \leq \epsilon \leq 1$$

$$\text{Say } \epsilon = 0.2$$

20% of the times choose the action randomly, not according to the Q-table.

$\epsilon \uparrow \Rightarrow$ agent explores more.

On-policy Monte Carlo Control

Find $\pi \approx \pi^*$

initialise $Q(s, a) \rightarrow$ random values,

loop many times,

① → simulate one episode using the Q -values in a ϵ -greedy method, $(s_0, a_1, r_1) (s_1, a_2, r_2) \dots$

② → set $G_1 = 0$

③ → update $Q(s, a)$ for the episode,
→ loop → $T-1, T-2, \dots, 0$

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

$$Q(s_t, a_t) = \underbrace{Q(s_t, a_t)}_{\text{new avg}} + \alpha \left(\underbrace{G_t - Q(s_t, a_t)}_{\substack{\downarrow \text{new value} \\ \downarrow \text{old avg}}} \right)$$

$$q_{\pi}(s, a) = \mathbb{E} \left[\underbrace{G_t}_{\text{new value}} \mid s_t = s, a_t = a \right]$$

