(✱). Contraction mapping Thm:-

$f : \mathbb{R}^n \to \mathbb{R}^m$

① 'f' which is contraction:-

$$|f(x) - f(y)| < (\alpha) \; |x - y|$$

↳ original.

$$0 \leq \alpha < 1$$

⇓

2. There is $\boxed{\text{a unique}}$ fixed Point.

$$\exists x : f(x) = x.$$

Method to compute fixed points:-

1. $\boxed{x_0 \in \mathbb{R}}$

2. Repeatedly apply 'f' on itself

$$x_0, \; f(x_0), \; f(f(x_0)), \; f(x_2) = f(f(f(x_0)))$$

$x_1$   $\downarrow x_2$   $x_3$

$$x_0, x_1, x_2, \cdots \longrightarrow \; (x^*)$$

⇓ satisfies

$$\boxed{f(x^*) = x^*}$$

Value Iteration :-

1. Identify a function.

$$V(s) = \boxed{\max_{a \in A}} \sum_{s' \in S} P(s'|s,a) \left[ r(s,a,s') + \gamma V(s') \right]$$

known    known    known.

✱ what is ⓥ & $\underline{f}$ ?

$$\left(\underline{V}\right) = \begin{bmatrix} V(s_1) \\ V(s_2) \\ \vdots \\ V(s_n) \end{bmatrix}$$

$$\begin{bmatrix} V(s_1) \\ V(s_2) \\ \vdots \\ V(s_n) \end{bmatrix} \quad f(\underline{v}) = \begin{bmatrix} \max\limits_{a} \sum\limits_{s' \in S} P(s'|s_1,a)\left(r(s_1,a,s') + \gamma V(s')\right) \\ \max\limits_{a} \sum\limits_{s' \in S} P(s'|s_2,a)\left(r(s_2,a,s') + \gamma V(s')\right) \\ \vdots \end{bmatrix}$$

$$\boxed{V = f(V).}$$

---

$\quad V_1, V_2 ; \rightarrow \|f(v_1) - f(v_2)\| \le \underline{\textcircled{\alpha}} \; \|v_1 - v_2\|. \quad \rightarrow \boxed{\text{Exercise}}$

$$\downarrow$$
$$\boxed{0 \le \alpha < 1}.$$

1. Start with any $V_0$.

2. $V_{i+1}(s) \leftarrow f(v_i)(s) = \max\limits_{a} \sum\limits_{s'} P(s'|s,a)\left(r(s,a,s') + \gamma \underline{V(s')}\right). \quad \hookrightarrow \boxed{\forall s \in S}$

3. $V_0, V_1, V_2, V_3, \cdots \rightarrow V^*$

$$\hookrightarrow \text{value iteration scheme.}$$

## chapter 5 :-

### Model - free techniques.

$\boxed{P(s'|s,a)} \rightarrow \boxed{\text{not known.}}$

Not known :- transition $\underline{\text{prob.}}$ , $\underline{\text{Structure of reward,}}$

given : samples / trajectories.

$$\rightarrow \left[ s_0, a_0, r_1, s_1, a_1, r_2, \cdots \right]$$

M.C prediction :-

$\hookrightarrow$ given a policy $\pi$

$\hookrightarrow$ Evaluate policy $\boxed{V^\pi}$

(assume determ. rewards)

$s, a \rightarrow r(s,a)$ is fixed.

$$V^\pi(s) = E\left[ \sum_{t=0}^{\alpha} \gamma^t r(s_t, a_t) \,\Big|\, s_0 = s, a_t \sim \pi \right].$$

$\boxed{(s_1, s_2, s_3 \cdots)}$

$\downarrow$

$s_0, a_0 \sim \Leftarrow \begin{bmatrix} \vdots \\ \vdots \end{bmatrix}$

Q: Sample average, Law of Large no's:-

$E[x]$; $\boxed{x_1, x_2, x_3, \cdots \underline{i.i.d.}}$

$H \to 0$
$\tau \to 1$.

$\boxed{\dfrac{\sum x_n}{n} \xrightarrow{a.s} E[x].}$

$\underset{\smile}{\textcircled{0}}, \underset{\smile}{\textcircled{1}} 1, 0, 1, 0, \cdots$

$X = \boxed{\sum\limits_{t=0}^{\infty} \gamma^t \cdot \gamma(s_t, a_t).} \Leftarrow t \atop \textcircled{T}$

$\nearrow$ Samples from this distribution

$E[x] = v^{\pi}(s)$.

$\boxed{\pi \atop v(s)}$

$\widehat{v^{\pi}(s)}$

1. generate an episode following $\pi$:

$\textcircled{$s_2$}$ $\quad\quad\quad\quad s_3 \atop \downarrow\downarrow$ $\quad \widehat{v^{\pi}(s)}$

1. $s_0 = s, a_0 \sim \pi, r_1', s_1', a_1' \sim \pi, r_2', \textcircled{$s_2$} \cdots \boxed{r_t'}$

step 2 $\qu\quad$ 2. $s_0 = s, a_0 \sim \pi, r_1^2, s_1^2, a_1^2, \cdots \quad\quad r_t^2$

$\vdots$

$1000 \cdot s_0^{1000} = s, - \cdots -$ $\quad\quad\quad\quad r_t^{1000}$.

step 2: $\boxed{\textcircled{$R_1$} = \dfrac{r_1' + \gamma r_2' + \gamma^2 r_3' + \cdots \gamma^{t-1} r_t'}{r_1^2 + \gamma r_2^2 + \cdots + \gamma^{t-1} r_t^2}}$

$\textcircled{$R_2$} =$

$\vdots$

$\textcircled{$R_{1000}$} = r_1^{1000} + \gamma r_2^{1000} + \cdots + \gamma^{t-1} r_t^{1000} \Big]$

$\boxed{v^{\pi}(s) \approx \dfrac{\sum\limits_i R_i}{i}}$ step 3

First visit:-

$(*)$ $\boxed{\boxed{s_1}\ a_1,\ 5,\ \boxed{s_3,\ a_?,\ 6}\ s_4,\ a_1,\ 3,\ \boxed{s_2}\ a,\ 4,\ \boxed{s_1,\ a_2,\ 8}}$ .

$V(s_1) \rightarrow\ \boxed{5 + \gamma \cdot 6 + 3 \gamma^2 + 4\gamma^3 + 8\gamma^4}\ ,\ 8$

$r(s_2) \leftarrow\ \textcircled{4} + \gamma 8$

$V(s_3) \leftarrow\ 6 + 3\gamma + 4\gamma^2 + 8\gamma^3$

$V(s_4) \leftarrow\ 3 + 4\gamma + 8\gamma^2$ .

**Every-visit :-**

$(*)$. Prediction : $\pi$ : $\underline{\underline{V^\pi}}$ .

$\downarrow$

Control :- Computing optimal policy $\underline{\underline{\pi^*}}$ using $\boxed{samples}$

**Monte Carlo Exploring starts :-**

$\boxed{V(s)} \longrightarrow \boxed{Q(s, a)} \leftarrow$

$\downarrow$

$\boxed{a \sim \pi}$

$S : \left(\begin{array}{c} Q(s, a_1) \\ Q(s, a_2) \\ \cdot \\ Q(s, a_n) \end{array}\right)$

$\boxed{\pi^*(s) = \arg\max_i Q(s, a_i)}$

$\boxed{\pi_0}$

$\downarrow$

$(*)$. 1. choose $s_0,\ a_0$ randomly.

2. generate an episode from $\underline{s_0, a_0}$ using $\boxed{\pi_0} \sim$ random policy.

$(*)$ $\boxed{for}$ each $\boxed{s\ \&\ action}$ $\boxed{(s, a)}$ :- $\boxed{randomly}$

$\begin{array}{l} 1.\ \boxed{s, a} \neq s_2 \cdots s_3 \cdots \\ 2.\ s, a, \quad-\quad-\quad \\ 3.\ s, a, \quad-\quad-\quad \end{array}$ Step 1.

step2:  $\begin{cases} R_1 \\ R_2 \\ \vdots \\ R_{1000} \end{cases}$

step 3:- $Q(S,a) \leftarrow \dfrac{R_1 + R_2 + \cdots R_n}{n}$

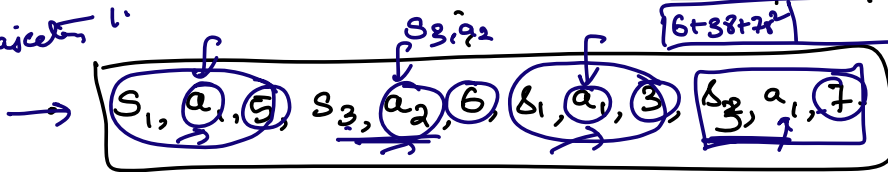$*$ $\pi_{i+1}(S) = \arg\max\limits_{a} Q(Sa)$ $\leftarrow$ improvement.

Q1:- Analogous to **Policy Iteration**

I. with same (trajectory)

$Q(S, a_1) \leftarrow$ | $R_1$ | $R_2$ | $4+8\gamma+6\gamma^2$ | $6..$ |

$S_1 a_2 \leftarrow$ | $0$ | | $\cdots$ |

$S_3 a_4 \leftarrow$ | $7$ | |

$\boxed{6+3\gamma+7\gamma^2}$

trajectory 1:

$\rightarrow$ $\boxed{S_1, a_1, 5, \ S_3, a_2, 6, \ S_1, a_1, 3, \ S_3, a_1, 7}$

$\boxed{\sum\limits_{t=0}^{\wedge} \gamma^t \, r_t}$

$R_1 = \boxed{5 + 6\gamma + 3\gamma^2 + 7\gamma^4}$; $R_2 = 3 + 7\gamma$.

$S_3, a_2 :-$ $6 + 3\gamma + 7\gamma^2$

$S_3, a_1 :-$ $(7)$

trajectory 2:- $S_1, a_1, 4, \ S_2, a_3, 8, \ S_1, a_1, 6.$

$\dfrac{4 + 8\gamma + 6\gamma^2}{6}$

$S, a_1 \leftarrow$ low rewards.

**online learning**

estimates of $\tilde{Q}_t(S, a). \ \forall S, a.$

on-policy first-visit MC control.

$\pi(S, a) \leftarrow \begin{cases} \arg\max\limits_{a} \tilde{Q}_t(S,a). \quad \text{w.p.} \\ \text{random action} \quad 1-\epsilon \, p. \end{cases}$