

(1). (*). UCB \rightarrow intelligent exploration.

✓ σ is low & σ is high. \rightarrow need for your intelligent exp.

⊕ σ is high & σ is low. ⊕

(2).

$$\pi P = \pi$$

$$P\pi = \pi$$

$$[x \ y] \begin{bmatrix} \frac{1}{3} & \frac{2}{3} \\ \frac{1}{2} & \frac{1}{2} \end{bmatrix} = [x \ y].$$

$$x + y = 1$$

$$\frac{x}{3} + \frac{y}{2} = x$$

$$\Rightarrow \frac{y}{2} = \frac{2x}{3} \Rightarrow y = \frac{4x}{3}$$

$$x + y = 1 \Rightarrow x + \frac{4x}{3} = 1 \Rightarrow \underline{\underline{x = 3/7}}$$

$$\underline{\underline{y = 4/7}}$$

(3A)

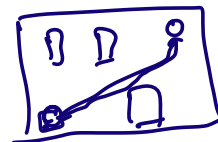
A : 1000 w.p. '1/1000' }
B : ②

$$E[r_A] = 1 ; E[r_B] = 2,$$

(4)

(c). Episodic ends after fixed no. of timesteps.

random



(5)

$$V_{\pi}(s) = \sum_a \pi(a|s) \cdot Q^{\pi}(s, a).$$

⊕s

$$V_{\pi}(s_1) = \frac{1}{3}(2+4+6) = 4.$$

$$V_{\pi}(s_2) = \frac{1}{3}(1+5+3) = 3.$$

$$\therefore V_{\pi}(s_1) + V_{\pi}(s_2) = \textcircled{7}.$$

Practice Questions :-

(1). (a). $\underbrace{E[r^i]} = (1-p_i)P_i + 0 \cdot (1-p_i)$
 $= \underline{(1-p_i)P_i}$

Arm 1: $\frac{3}{4} \cdot \frac{1}{4}$

2: $\frac{1}{3} \cdot \frac{2}{3}$

3: $\frac{1}{4} \cdot \frac{1}{4}$ ✓

4.
5.

(b). $\boxed{f(p)} = (1-p)p$ ✓

$1-2p=0 \Rightarrow \boxed{p=\frac{1}{2}}$

$-2 < 0$ ✓

(2).

Policy Iteration :-

↑ stochastic.

$s, a \rightarrow r$ is fixed

1. Evaluating a given policy (π):

$$\boxed{V^{\pi}(s)} = \sum_a \pi(a|s) \cdot \sum_{s'} p(s'|s,a) (r(s,a) + \gamma V^{\pi}(s'))$$

↳ system of equations

(*) Contraction Mapping thm :-

• f is a contraction :

↓

• \textcircled{f} has a fixed point :
unique.

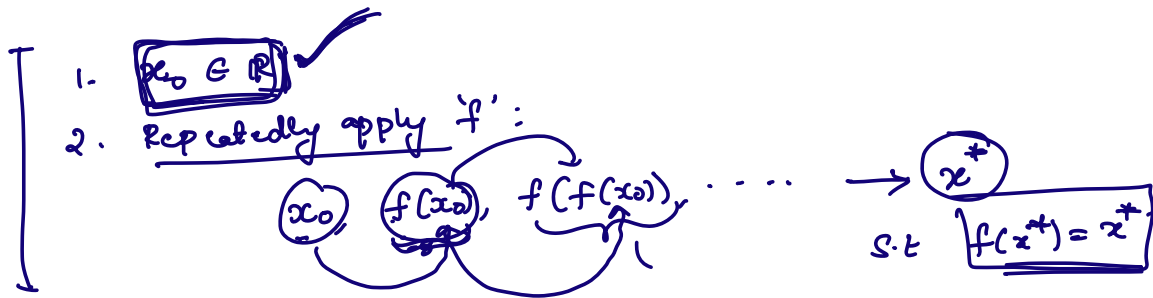
\textcircled{x} : $\underbrace{f(x)} = \underbrace{x}$

$$\|f(x) - f(y)\| \leq \alpha \|x - y\|$$

$\boxed{0 < \alpha < 1}$

$\begin{matrix} x & \longrightarrow & y \\ \downarrow & & \downarrow \\ f(x) & \longrightarrow & f(y) \end{matrix}$

how to find fixed point :-



$$f(V) = \sum_a \pi(a|s) \sum_{s'} P(s'|s,a) \cdot (r(s,a) + \gamma V(s'))$$

$V = f(V) \rightarrow$ fixed point of this function

(*) Show f is contraction. $\sqrt{\|f(V_1) - f(V_2)\|} \leq \alpha \|V_1 - V_2\|$
 \downarrow
 $0 < \alpha < 1$

Iterative policy evaluation:-

1. Initialising $V(s) \equiv 0$

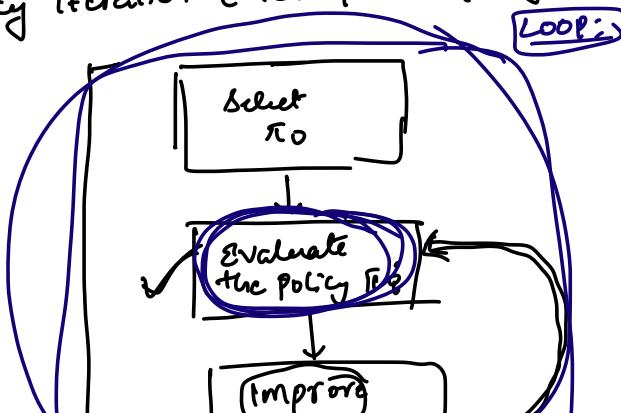
2. Loop:- \rightarrow 3 step

Loop for each state:-

$$V'(s) \leftarrow \sum_a \pi(a|s) \sum_{s'} P(s'|s,a) (r(s,a) + \gamma V(s'))$$

3. $V^{\pi}(s) \leftarrow \text{result}$

Policy iteration (for optimal policy):-



(*) $\{ \frac{1}{n} \} \rightarrow 0$

$1, \frac{1}{2}, \frac{1}{3}, \dots$
 $\frac{1}{2}, \frac{1}{4}, \dots$

$S, (A)$ $1, 2, ?$
 A, A, A, \dots

deterministic? (A^S)

$\Pi: S \rightarrow A$

The Policy Mix

(ii) is the optimal

TICS
(new policy)

$$v^{\pi'} \geq v^{\pi}$$

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

(2 A).

$V_o(s) = 0, \forall s.$

Step 1:

$$\pi : [1/4, 1/2, 1/2, 1/4]$$

$$\underline{V_1(t)} := \frac{1}{4} [-1 + 8V_0(1)] + \frac{1}{4} [-1 + 8V_0(2)] + \frac{1}{4} [-1 + 8V_0(3)] + \frac{1}{4} [-1 + 8V_0(5)].$$

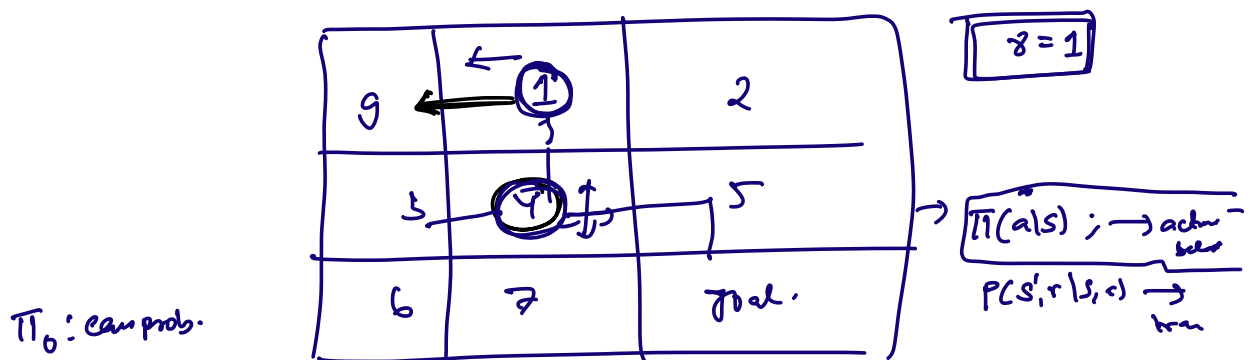
Step 2:

Step 3:-

(b). Policy improvement :-

$$\boxed{V_1(1) = V_1(2) = V_1(3) = V_1(5) = -1.25.}$$

$$V_1(4) = -1. ; V_1(6) = V_1(7) = -1.5$$



$$(*) \pi_1 := \arg \max_{\pi} \left\{ \sum_{s', r} P(r, s' | s, a) \left(\underbrace{r}_{\text{description } P.} + \gamma \underbrace{V^{\pi_0}(s')}_{\text{best.}} \right) \right\}.$$

$$\arg \max_{\{U, D, L, R\}} \left\{ \begin{array}{l} -2 + 1(-1.25), \\ -1 + 1(-1), \\ -1 + 1(0), \\ \boxed{-1 + 1(-1.5)}. \end{array} \right\} \rightarrow \text{best.}$$

$$\pi_1(4): \arg \max_{\{U, D, L, R\}} \left\{ \begin{array}{l} -1 + 1(-1.25), \\ -1 + 1(-1.25), \\ -2 + 1(-1.25), \\ -1 + 1(-1.25) \end{array} \right\}.$$

(3). Status:- { high, low }
 ↓
 action:- { (wait, search), (recharge, search, wait) }.

$$V^*(s) = \max_a \sum_{s', r} P(s', r | s, a) [r + \gamma V^*(s')].$$

$$V^*(\text{high}) = \max_{\{\text{wait}, \underline{\text{search}}\}} \left[\begin{array}{l} r_{\text{wait}} + \gamma V^*(\text{high}), \\ \alpha [r_{\text{search}} + \gamma V^*(\text{high})] + \\ \underline{(1-\alpha)} [r_{\text{search}} + \gamma V^*(\text{low})]. \end{array} \right]$$

$$V^*(\text{low}) = \max_{\{\text{recharge}, \underline{\text{wait}}, \underline{\text{search}}\}} \left[\begin{array}{l} 0 + \gamma V^*(\text{high}), \\ r_{\text{wait}} + \gamma V^*(\text{low}), \\ \beta [r_{\text{search}} + \gamma V^*(\text{low})] + \\ \underline{(1-\beta)} [-3 + \gamma V^*(\text{high})]. \end{array} \right]$$

Bellman eq. for optimal policy.

(b). $\pi = (\text{search}, \text{recharge})$

V^π (Bellman eq. for given policy).

$$\left. \begin{array}{l} V^\pi(\text{high}) = -10 + 3 \cdot V^\pi(\text{low}). \\ V^\pi(\text{low}) = 5 + 4 \cdot V^\pi(\text{high}). \end{array} \right\}$$