Recap :- Bandits → e-greedy & UCB.

## MDP :-

- Agent , Environment.

- Sequential Decision Making.
  - trajectory :

$$(s_0, a_0, r_1, s_1, a_2, r_3, s_2, \dots).$$
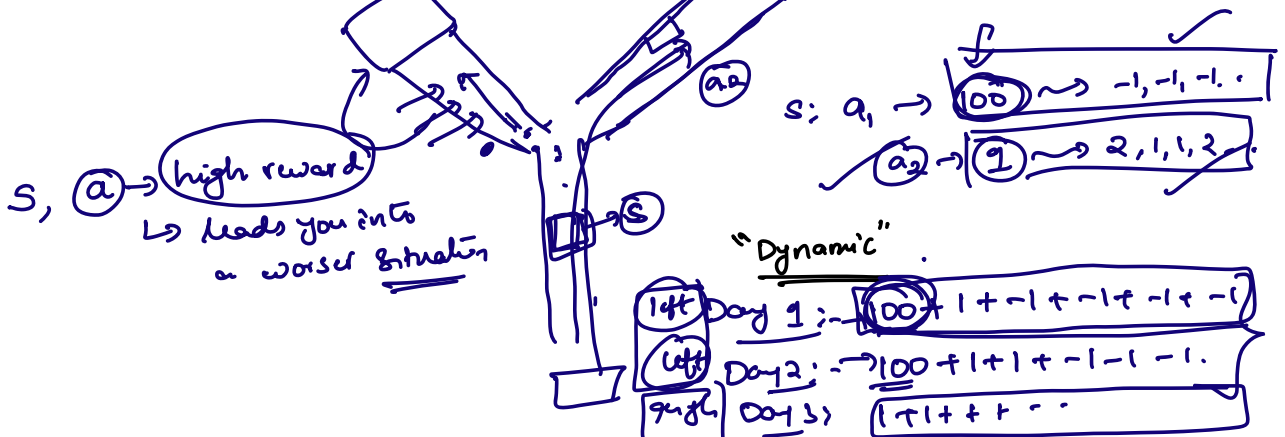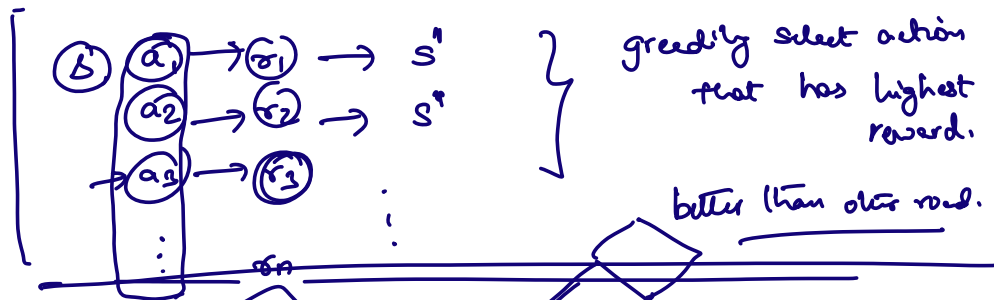
- (Bandit) vs MDP :-

$$s, (a_1, a_2, \dots a_n), r.$$

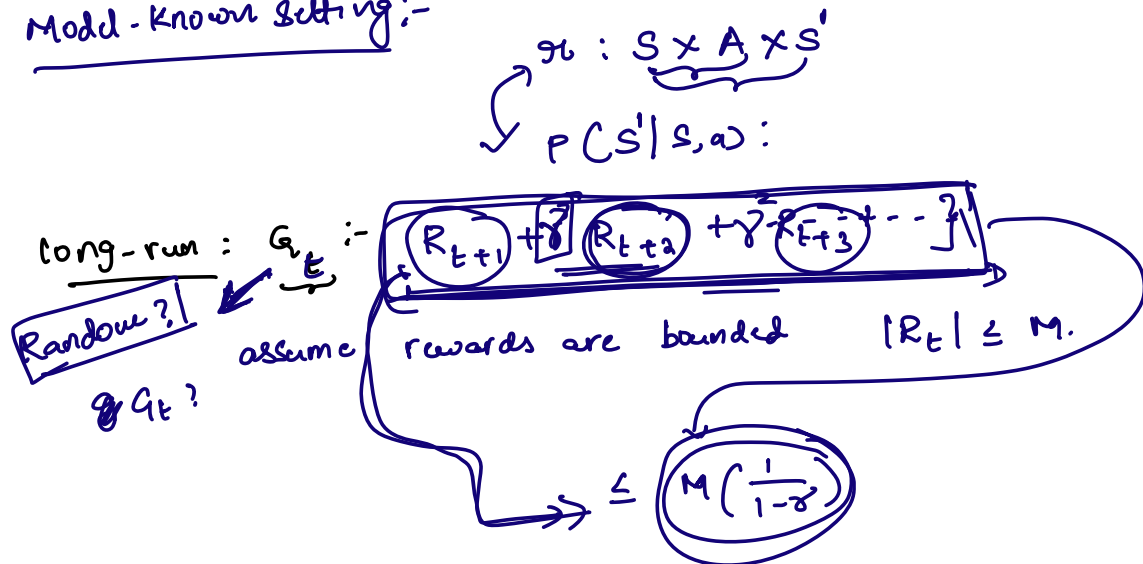$$\boxed{P(s|s,a_i) = 1} \qquad \underline{\text{best arm.}}$$

- MDP:- Sequential decision Making
  ⑤ Find what is the best action/decision to make.

  " long-run goal "



greedily select action that has highest reward.

better than other road.

$S, \text{ⓐ} \rightarrow$ high reward
  ↳ leads you into a worser situation

$s; a_1 \rightarrow \boxed{100} \rightsquigarrow -1, -1, -1 \dots$
$\qquad a_2 \rightarrow \boxed{1} \rightsquigarrow 2, 1, 1, 3 \dots$

"Dynamic"

Day 1 :- $100 + 1 + -1 + -1 + -1 + -1$
Day 2 :- $100 + 1 + 1 + -1 - 1 - 1.$
Day 3 :- $1 + 1 + + + \dots$

obj:- Select actions that has higher long-run rewards.

---

## Model-Known Setting:-

$$\pi : S \times A \times S'$$

$$P(S'|S,a):$$

long-run : $G_t$ :- $[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \cdots ]$

Random?!

& $g_t$?

assume rewards are bounded    $|R_t| \leq M.$

$$\leq M\left(\frac{1}{1-\gamma}\right)$$

## Machine Replacement:-

(*).  $\boxed{r : S \times A \longrightarrow \mathbb{R}}$

Policy:   $\pi : \textcircled{S} \longrightarrow A$  ;    $\pi : S \times A \longrightarrow \mathbb{R}$  s.t $\sum\limits_{a} \pi(s,a) = 1.$

deterministic

stochastic

$$\pi(s, a_1) = 0.6$$
$$\pi(s, a_2) = 0.4.$$

$a_0$ :- $s_0, a_0, \textcircled{s_1}, a_1, \cdots$

$\rightarrow P(\cdot | s_0, a_0)$  ← P.T.M.

$\textcircled{E}$ $[r(s_0, a_0) + \gamma r(\textcircled{s_1}, a_1) + \gamma^2 r(\textcircled{s_2}, a_3) + \cdots]$

$\dfrac{s_1, s_2, \cdots}{a}$

---

fix actions

random new how wind

goal.

stall = $\begin{bmatrix} \text{current} \\ \text{distance from nearest} \\ \text{obstacle} \\ \text{distance from goal} \end{bmatrix}$

actions:- L, R, B, F

$1^{st}$ trajectory, $S_0$, left, $S_1$, left,
$\gamma$

$(\ast)$. $\boxed{V^{\pi}(s)} = E\left[\sum_{t=0}^{\infty} \gamma^t \, r(\underline{S_t, a_t}) \mid \underline{S_0 = s}, a_t = \pi(S_t)\right].$

$(S_1, S_2, S_3 \cdots)$

total discount reward  |  Starting from state $a$  |  following the policy $\pi$

---

Assume there is only one policy.

$\boxed{\text{Single action}}$ in every state.

$\boxed{A = \{a\}.}$

$\boxed{r(s) = r(s,a).}$

---

$V(s) = \underset{(S_1, S_2, \cdots)}{E}\left[\sum_{t=0}^{\infty} \gamma^t \, r(S_t) \mid S_0 = s\right]$

$\longrightarrow \gamma^0 r(s).$

$\boxed{r(s)}$

$= \sum_{(S_1, S_2, \cdots)} \Pr\{S_1 = \delta_1, S_2 = \delta_2, \{s=s\}\} \sum_{t=0}^{\infty} \gamma^t \, r(S_t).$

$= \underline{r(s)} + \sum_{(S_1, S_2, \cdots)} \Pr\{S_1 = \delta_1, S_2 = \delta_2, \cdots \mid S_0 = s\} \cdot (\gamma r(S_1) + \gamma^2 r(S_2) + \gamma^3 r(S_3) + \cdots)$

$= r(s) + \gamma \sum_{(\delta_1, \delta_2, \cdots)} \Pr\{S_1 = \delta_1, S_2 = \delta_2, \cdots \mid S_0 = s\} \left[r(S_1) + \gamma \, r(S_2) + \cdots\right]$

---

$\Pr\{\underset{A}{S_1 = \delta_1}, \underset{B}{S_2 = \delta_2}, \cdots \mid \underset{c}{S_0 = s}\} \overset{?}{=}$

$\Pr\{S_2 = \delta_2, S_3 = \delta_3, \cdots \mid \boxed{S_1 = \delta_1}, \cancel{S_0 = s}\}.$ ← Markov property

$\Pr\{S_1 = \delta_1 \mid S_0 = s\}.$

$\Pr\{A, B \mid c\} = \Pr\{A \mid c\} \ast \Pr\{B \mid A, c\}.$

$= r(s) + \gamma \sum_{(\delta_1, \delta_2, \cdots)} \boxed{\Pr\{S_1 = \delta_1 \mid S_0 = s\} \ast \Pr\{S_2 = \delta_2, S_3 = \delta_3 \mid S_1 = \delta_1.\}}$

$\left[r(S_1) + \gamma r(S_2) + \cdots\right]$

$$= r(s) + \gamma \sum_{s_1} Pr\{S_1 = s_1 | S_0 = s\} \cdot \quad \sum_{(s_2, s_3, \ldots)} Pr\{S_2 = s_2, S_3 = s_3, \ldots S_\infty = 1, 1\} \cdot [r(s_1) + \gamma r(s_2) + \gamma^2 r(s_3) + \ldots]$$

all the states that can occur at time $i$.

$$= \boxed{r(s) + \gamma \sum_{s'} Pr\{S = s' | S_0 = s\} \cdot V(s')}$$

$\pi$ : determistic.

$$V^\pi(s) = r(s, \pi(s)) + \gamma \sum_{s'} Pr(s' | s, \pi(s)) \, V^\pi(s'). \qquad \rightarrow \text{Bellman equation}$$

$$\boxed{Q^\pi(s, a)} = r(s, a) + \gamma \sum_{s'} Pr(s' | s, a) \cdot V^\pi(s''). \qquad \rightarrow \text{Q-Bellman eqn.}$$

$\pi$ : stochastic :

$$V^\pi(s) = \boxed{\sum_{a \in A} \pi(a|s)} \left[ r(s, a) + \gamma \sum_{s'} Pr(s' | s, a) \cdot V^\pi(s') \right]$$

$$\text{(★★★)} \quad Q^\pi(s, a) = r(s, a) + \gamma \sum_{s'} Pr(s' | s, a) \cdot V^\pi(s'). \longleftarrow \quad ①$$

$$\pi(s) \rightarrow \begin{array}{l} a_1 \rightarrow \frac{1}{2} \\ a_2 \rightarrow \frac{1}{2} \end{array}$$

$$\boxed{V^\pi(s)} = \boxed{\sum_{a \in A} \pi(a|s) \cdot Q^\pi(s, a).}$$

$$\boxed{\frac{1}{2} \cdot Q^\pi(s, a_1) + \frac{1}{2} Q^\pi(s, a_2).}$$

$$\boxed{V^\pi(s) = \sum_{a \in A} \pi(a|s) \cdot Q^\pi(s, a).} \qquad \longrightarrow \quad ②$$

w.r.t given policy $\pi$.

$$Q^{\pi}(s,a) = r(s,a) + \gamma \sum_{s',a'} Pr\{s'|s,a) \pi(a'|s') \cdot Q^{\pi}(s',a') \longrightarrow \textcircled{3}$$

$\longrightarrow$ Bellman eqn for optimal policy

$\longrightarrow$ Policy Iteration / V.I.

Solve proactive $\Uparrow$