

Question 3)

Ex 3.15)

NO, signs are not important ~~only~~
Interval between them is important

$r_E \rightarrow$ reward of hitting edge
 $r_A \rightarrow$ reward for leaving box A or B.
 $r_B \rightarrow$ " " " B
 $r_0 \rightarrow$ " " other action's

~~the~~
Suitable order

$$r_E < r_0 < r_A, r_B$$

$$V_k = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n}$$

$$V_k = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n V_{t+n}$$

$$V_c = V_k + c [1 + \gamma + \gamma^2 + \dots + \gamma^{n-1}]$$

$$V_c = c [1 + \gamma + \gamma^2 + \dots + \gamma^{n-1}]$$

$$V_c = c + c\gamma + c\gamma^2 + \dots + c\gamma^{n-1}$$

$$- \gamma V_c = -c\gamma - c\gamma^2 - \dots - c\gamma^n$$

$$V_c - \gamma V_c = c(1 - \gamma^n)$$

$$V_c = \frac{c(1 - \gamma^n)}{1 - \gamma} = \frac{c}{1 - \gamma}$$

Since
 $\lim_{n \rightarrow \infty} \gamma^n = 0$
 $|V_c|$

Ex 3.16 >

It will change the Maze running
Because in episodic task n is
finite.

from equation's of 3.15

$$\cancel{V_*^{\text{original}}} V_*^{\text{modified}} = V_*^{\text{original}} + \frac{c(y^n - 1)}{1-y}$$

so for different state's n will
have different Value's.
Highest for starting state and
lowest for least visited state.

Q 5)

$$V_{\pi}(s) = \sum_a \pi(a|s) \left(\sum_{\mathcal{R}} \sum_{s'} p(s', \mathcal{R} | s, a) [\mathcal{R} + \gamma V_{\pi}(s')] \right)$$

$$= \sum_a \cancel{q} \sum_a \pi(a|s) q_{\cdot}(s, a)$$

Since $q(s, a) = \sum_{\mathcal{R}} \sum_{s'} p(s', \mathcal{R} | s, a) [\mathcal{R} + \gamma V_{\pi}(s')]$

Ans 8 >

Yes, R_{t+2} depends of S_t, A_t

Let S, R, A be the set of state, Reward action.

S'_1, S'_2, \dots, S'_n be represent state's after jumping to next state from initial state (S) by taking action A, a .

$$\pi(R_{t+2} | S_t = s, A = a)$$

$\pi(\cdot)$ \rightarrow this represent PMF of R_{t+2}

Below expression represent probability being in S'_i after state S after taking action a .

$$p(S'_i) = \sum_r p(S'_i, r | s, a)$$

$p(\cdot)$ is PMF. state is random variable

$$E(S'_i) = \sum \sum \pi(a | S'_i) \left(\sum_{S'} \sum_{R} r p(S', R | S'_i, a) \right)$$

$E(\cdot)$ is ~~PMF~~ Expected reward in state S'_i

$$E[R_{t+2} | S_t = s, A_t = a] = \sum_{i=1}^N p(S'_i) E(S'_i)$$

Since $p(S'_i)$ depends on S_t and A_t so, R_{t+2} also depend on S_t and A_t .

'N' is total number of states

Ans 9)

$$E[R_{t+2} | S_t = s, A_t = a] =$$

$$\sum_{i=1}^N p(s'_i) E(s'_i)$$

from Ans 8.

Ans 10)

$\pi(a|s)$ at $A, s \in S$ is given
 $p(s', r | s, a)$ is also given.

$$V_{\pi}(s) = E(G_t | S_t = s)$$

$$= \sum_a \pi(a|s) E(G_t | S_t = s, A_t = a)$$

Solving for $E[G_t | S_t = s, A_t = a]$

$$= E[R_{t+1} + \gamma G_{t+1} | S_t = s, A_t = a]$$

$$= E[R_{t+1} | S_t = s, A_t = a] + \gamma E[G_{t+1} | S_t = s, A_t = a]$$

$$= \sum_{s'} \sum_r r p(s', r | s, a) + \gamma \sum_{s'} \sum_r p(s', r | s, a) V_{\pi}(s')$$

— (ii)

Replace $E[G_t | S_t = s, A_t = a]$ by ~~from~~ eq (ii)

$$V_{\pi}(s) = \sum_a \pi(a|s) \left(\sum_{s'} \sum_{a'} p(s', a' | s, a) [r + \gamma V_{\pi}(s')] \right)$$

Ans 11)

$$= 2$$

$$= 0.5(-1) = -0.5$$

$$= (0.5)^2(10) = 2.5$$

$$= (0.5)^3(-3) = -0.375$$

Let $E(G_t)$ be x

$$x = c + \gamma c + \gamma^2 c + \dots + \gamma^{n-1} c \quad \text{--- (i)}$$

$$x\gamma = \gamma c + \gamma^2 c + \dots + \gamma^n c \quad \text{--- (ii)}$$

subtract (ii) from (i)

$$x(1-\gamma) = c - \gamma^n c$$

$$x = E[G_t] = \frac{c(1-\gamma^n)}{1-\gamma}$$

Since $|y| < 1$ and $n \rightarrow \infty$ so, $\lim_{n \rightarrow \infty} \gamma^n \rightarrow 0$

$$E[G_{\infty}] = \frac{c}{1-\gamma}$$

Ans 1.2) We have $V_*(s)$ $s \in S$.

for any state s in S .

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

Below

take action 'a' which maximises above expression for a given state.

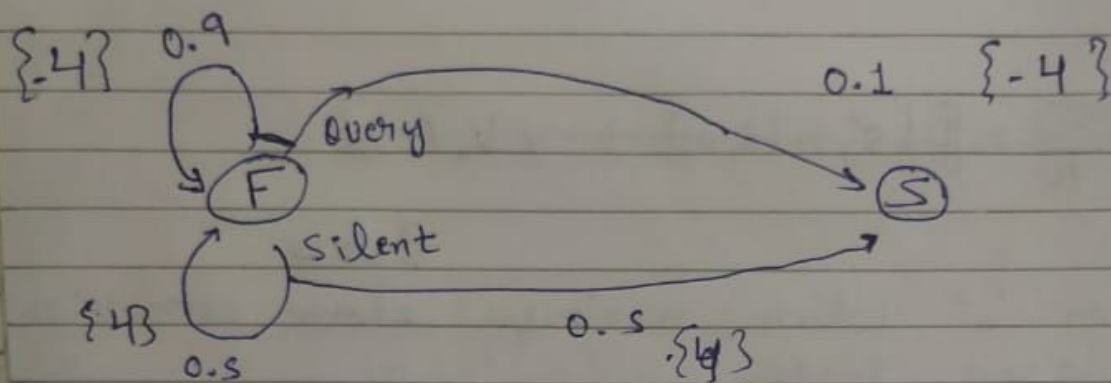
if new $V_*(s)$ is greater than old $V_*(s)$

then replace old with new one.

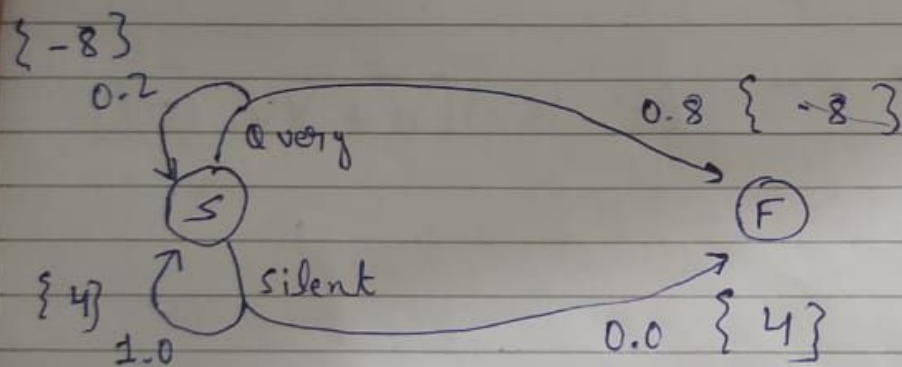
$$\arg \max_a \sum_{s'} \sum_{r} p(s', r | s, a) [r + \gamma V_*(s')]$$

Q 13

states = { fresh, stale }



{ I } represent reward reward



cost is converted into reward by changing sign.