



Task 1.

Action State	up	down	left	right
E	0	0	0	0
T	0	0	0	0
S	0	0	0	0
B	0	0	0	0

$$Q=0.1; \gamma=1$$

$$R_E=-1; R_T=+100; R_B=-0.5; R_S=-1$$

$$Q_{k+1}^{\pi}(X_t, u_t) = (1-\alpha) \cdot Q_k^{\pi}(X_t, u_t) + \alpha \cdot (R_{t+1} + \gamma \max_u Q_k^{\pi}(X_{t+1}, u))$$

By neglecting epsilon-greedy strategy, which action should be chosen? $\max Q$ or random u ?

Q-Table of Iteration "0"

Step 2: Choose an action = e.g. $Q(S, up)$ with $X_1=S$

Step 3: Perform an Action: $X_2=E$

all Q-values are zero
Randomly select.

Step 4: Measure rewards: $R_E=-1$

Step 5: Update Q-Table: $Q_2^{\pi}(S, up) = (1-0.1) \times 0 + 0.1 \times (-1 + 1 \times 0)$

$$= -0.1$$

State E is now the current state.

Action State	up	down	left	right
E	0	0	0	0
T	0	0	0	0
S	-0.1	0	0	0
B	0	0	0	0

Q-Table of Iteration "1"

Task 2

Action State	up	down	left	right
E	0	-0.1	-0.1	0
T	0	0	0	0
S	-0.1	-0.1	-0.1	7.91
B	40.95	0	-0.1	0

Q-Table of Iteration "0"

$$\alpha = 0.1; \gamma = 1$$

$$r_E = -1; r_T = +100; r_B = -0.5; r_S = -1$$

$$Q_{k+1}^{\pi}(X_t, u_t) = (1 - \alpha) \cdot Q_k^{\pi}(X_t, u_t) + \alpha (r_{t+1} + \gamma \max_n Q_k^{\pi}(X_{t+1}, u))$$

Step 2: Choose an action = e.p. $Q(S, \text{right})$ with $X_1 = S$

Step 3: Perform an Action: $X_2 = B$

choose the largest Q-value in State S

Step 4: Measure rewards: $r_B = -0.5$

Step 5: Update Q-Table: $Q_2^{\pi}(S, \text{up}) = (1 - 0.1) \times 7.91 + 0.1 \times (-0.5 + 1 \times 40.95)$

$$= 11.764$$

State B is now the current State.

⇒

Action State	up	down	left	right
E	0	-0.1	-0.1	0
T	0	0	0	0
S	-0.1	-0.1	-0.1	11.764
B	40.95	0	-0.1	0

Q-Table of Iteration "1"

Task 3 : ϵ -Greedy Strategy

$$\text{Action} = \begin{cases} \max_u Q(x, u), & R > \epsilon \Rightarrow R \text{ is uniform random value in } [0, 1] \\ \text{Random } u, & R \leq \epsilon \Rightarrow \epsilon(\text{episode}) = 0.9999^{\text{episode}} \end{cases}$$

- At the beginning, ϵ (epsilon) is close to 1, and the random number "R" is likely to be less than " ϵ ", "Action" tends to be randomly selected.
- When the episode is large enough, the " ϵ " (epsilon) becomes very small, and the random number "R" is likely to be greater than " ϵ ". The "Action" tends to be selected as the action with the largest Q-value in the current state.
- Q-Table: state S: $Q(S, \text{left})$ will be ^{the} largest, $Q(S, \text{up})$ will be larger than the rest two.
state E: $Q(E, \text{right})$ will be ^{the} largest.
state B: $Q(B, \text{up})$ will be the largest.
state T = Target location, all Q-Value = 0

Action \ State	up	down	left	right
E	-X	-X	-X	100
T	0	0	0	0
S	99	-X	-X	99.5
B	100	-X	-X	-X

Task 4

o When the episode is large enough, every Q -value update path will not pass through the square E , what means, the Q -value corresponding to state E will no longer be updated.









