

## # Algo

- i) Initialize the policy iteration, environment, discount factor ( $\gamma$ ) and initial policy over actions for each state.
- ii) Policy Evaluation:
  - For a given state, calculate the value of each action.
  - Update the value of current state using the calculated action values.
  - Calculate the maximum change in the value function.
  - Repeat if the maximum change is above the threshold.
- iii) Policy Improvement:
  - For a given state, calculate the value of each action.
  - Update the policy to choose the action with the highest value.
- iv) Repeat step 2 & Step 3 until convergence.
- v) Print result (Optimal policy).

## Exp 3:

Algo: i) Initialize the MDP environment, discount factor ( $\gamma$ ) & the value function.

ii) Repeat until convergence

- For a given state, calculate the value of each action.

$$Q_H(s, a) = H(a|s) * (R + \gamma V_H(s))$$

- Update the value function for each state using optimal action value.

$$V_H(s) = \max_a [Q_H(s, a)]$$

- Exit if maximum change involve function is below threshold value.

iii) For each state, select the action that maximise the action value  $H(s) = \arg \max_a (Q_H(s, a))$

iv) Print the result (Optimal policy).

## Exp 4:

same as Exp 2

v) print result (Optimal policy, value function).

## Exp 5: Algorithm:-

i) Initialize the Q-learning environment, learning rate, discount factor, epsilon and the Q-function.

ii) Repeat for all the episodes till convergence.

- Decide the trade off & choose appropriate action.
- Use Bellman Equation to update the Q-function.

iii) Find the optimal policy which has the maximum Q-value.

iv) Print the result (Q-function & optimal policy).

Exp 6: Bellman equation :-

Algorithm: same as MDP

iv. & print (Value function & optimal policy)

Exp 7: Monte Carlo :-

Algo:-

i. & Initialize Monte Carlo environment, gamma, epsilon & Q-function

ii. & Repeat for all the episodes:

- Decide the stochastic and generate an episode until it terminates

- Update the Q-function using the episode history & observed returns

iii. & Find the optimal policy which has maximum Q-value

iv. & Print the result (Q-function + Optimal Policy)

$$Q(s, a) \leftarrow Q(s, a) + \underset{\substack{\uparrow \\ \text{learning} \\ \text{rate}}}{\alpha} \cdot (\underset{\substack{\uparrow \\ \text{Discount} \\ \text{factor}}}{\gamma} \cdot (\underset{\substack{\uparrow \\ \text{next} \\ \text{action}}}{\max_{a'} Q(\overset{\substack{\uparrow \\ \text{next state}}}{s'}, \overset{\substack{\downarrow \\ \text{next state}}}{a'})}) - \underset{\substack{\uparrow \\ \text{current} \\ \text{action}}}{Q(\overset{\substack{\downarrow \\ \text{current state}}}{s}, \overset{\substack{\downarrow \\ \text{current state}}}{a})}))$$