

(18)

RL framework: The solution



once a problem is properly defined, the solution can be thought of



solution to a RL problem can be considered as a series of actions that need to be learned (or taken) by the agent in order to maximize the reward signal.



these actions depend on the state of the task and hence being able to come up with the best possible action for a given state, can be considered to be a solution to the RL problem.



motivates the idea of policy i.e. mapping b/w the states and actions.

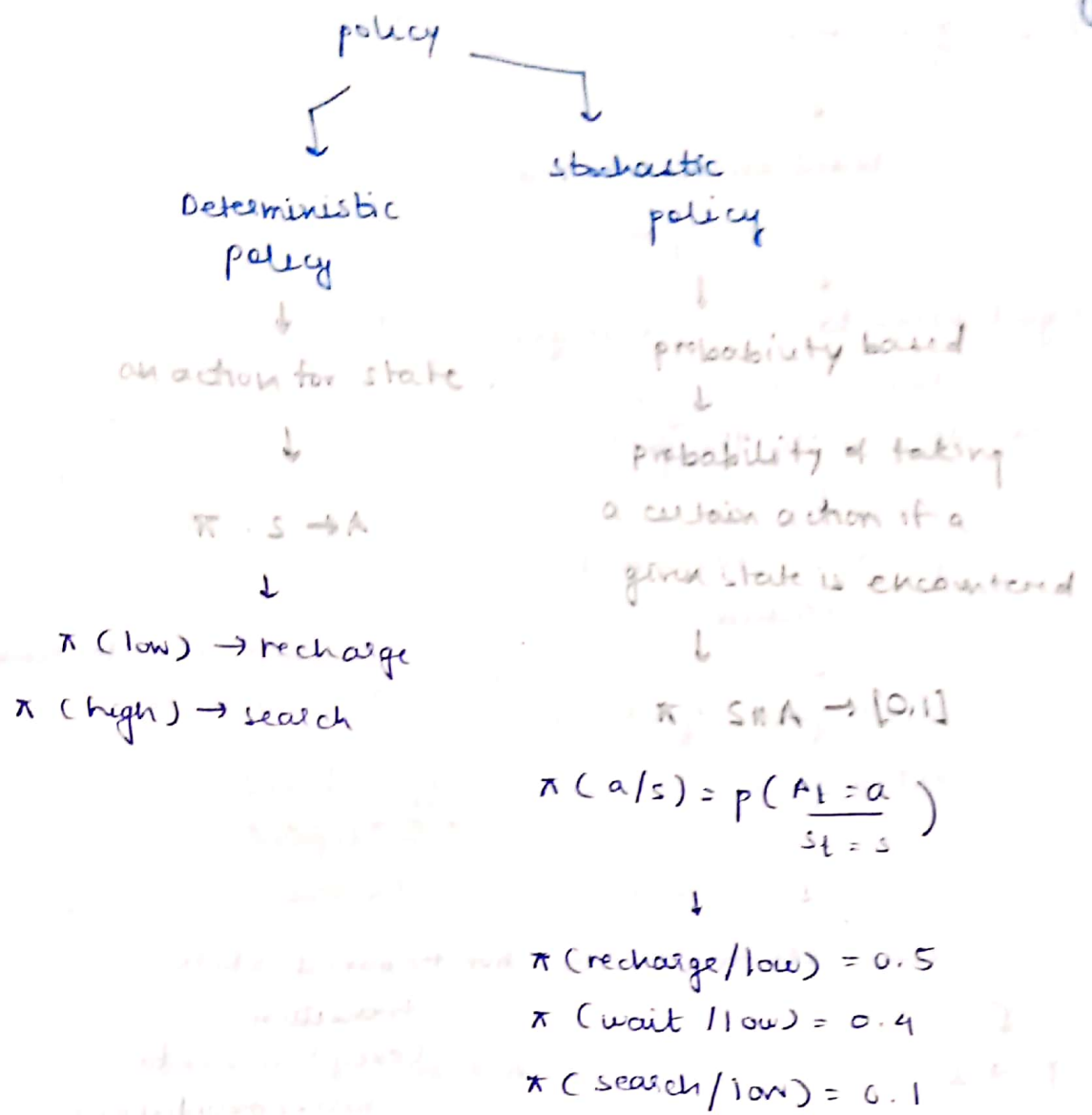


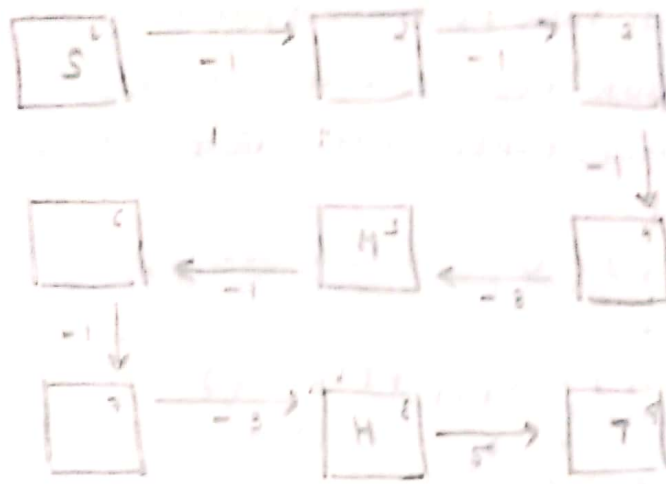
obviously, there will be many possible, but not all will be optimal.

→ finding this is the key requirement



⇒ optimal policy for a RL problem is considered as its solution.





Approximation only.
 how do we know if it's comp? better? $CR = -6$
 policy \rightarrow action value (to input policy)
 state value function
 (to compare with other policies)

state value function

The process is repeated for the cell right adjacent to S, and then get $CR = -5$, and all cells soon.

↓
 for terminal state, $R = 0$, since if agent starts there, the episode ends straight away.

↓
 Accordingly, we have a number associated to each state, this number is called the output of the state value function.

↓
 state value function yields the reward expected if the agent were to start at that position & follow the given policy.

↓

$$v_{\pi}(s) = E_{\pi} \left[\sum_{t=0}^{\infty} \gamma^t r_t \mid s_t = s \right]$$
 Mathematical notation for state value function

↓
 ** for each state s , it yields the expected discounted return if the agent starts in state s and then uses the policy to choose its actions for all time steps.

(22)

State value function have a recursive property where the state value function output is equal to the immediate reward added to the discounted reward of the next state ($\gamma = 1$)

↓

$$V_{\pi}(s_9) = 0 \quad (\text{since terminal})$$

↓

$$V_{\pi}(s_8) = R_{t+1} + \gamma V_{\pi}(s_9) \\ = 5$$

↓

$$V_{\pi}(s_7) = R_{t+1} + \gamma V_{\pi}(s_8) \\ = -3 + 1(5) = 2$$

↓

$$V_{\pi}(s_6) = R_{t+1} + \gamma V_{\pi}(s_7) \\ = 2 + 1(2) = 4$$

↓

$$V_{\pi}(s_5) = R_{t+1} + \gamma V_{\pi}(s_6) \\ = -1 + 1(4) = 3$$

↓

$$V_{\pi}(s) = E_{\pi} \left[\frac{R_{t+1} + \gamma V_{\pi}(s+1)}{s_t = s} \right]$$

↓

Bellman equation

↓

since the next state on taking a certain action depends on one step dynamics and our state

value is dependent on the next state,

then it should be accounted in the expected reward and the corresponding rewards

* Also the action taken at a given state, which eventually decides the next possible states, may vary depending on the type of policy which can be either deterministic or stochastic.

↓
↳ so that needs to be taken into account as well.

i) deterministic → one action per state.

future reward → depends on the state, action pair, and that is also dependent on the type of policy.

now the probability factor comes in for the resulting state, reward for that action, which is again given by the one step dynamics of the environment

$V_{\pi}(s)$ ↓
state value for current values as weighted sum of possible next states given the current state.

$$V_{\pi}(s) = \sum_{\substack{s' \in S^+, \\ r \in R}} p\left(\frac{s', r}{s, \pi(s)}\right) (r + \gamma V_{\pi}(s'))$$

↓
one step dynamics.
↓
probability of s' with corresponding r when current state is s and action $a \rightarrow \pi(s)$, which is deterministic

ii) stochastic → probability based action, given a state, adds another probability factor for selection of action.

↓
once action is selected, the further process is same.

(weighted sum) ↑

$$V_{\pi}(s) = \sum_{\substack{s' \in S^+, \\ r \in R, \\ a \in A(s)}} \pi(a/s) p\left(\frac{s', r}{s, \pi(s)}\right) (r + \gamma V_{\pi}(s'))$$

↓ prob of action ↓ prob of new state for the action

↑
reward for (s, a)

②④ → Comparing policies

↓
given two policies π_1 and π_2 , we can say that

π_1 is better than π_2 , if $V_{\pi_1}(s) > V_{\pi_2}(s)$

for all $s \in S$.

↓
in some situations, it might be hard to compare policies.

↓
But an optimal policy will always be available for a problem, denoted by π^*

Action value function

↓
value of taking action a in state s under policy π

↓
if action a is taken at state s , then what is the reward, if the policy π is followed, for the next state.

↓
a value for each possible action in a state s .

↓

$$q_{\pi}(s, a) = E_{\pi} \left[\frac{G_t}{s_t, a_t} \right]$$

↓

$s \quad a$

By choosing the best action value per state, we move from q_{π} to π^*

next problem is,

how do we determine this ~~best~~ policy from bad policy

↓
this is done through action value function