

# Project Tasks.

Q1 Derive Bellman optimal equations for value & action value functions.

$$\begin{aligned} V_*(s) &= \max_{a \in A} q_{\pi_*}(s, a) = \max_{a \in A} E_{\pi_*}[G_t | S_t = s, A_t = a] \\ &= \max_{a \in A} E_{\pi_*}[R_{t+1} + \gamma V_{t+1} | S_t = s, A_t = a] \\ &= \max_{a \in A} E_{\pi_*}[R_{t+1} + \gamma V_*(s_{t+1}) | S_t = s, A_t = a] \end{aligned}$$

$$V_*(s) = \max_{a \in A} \sum_{s', r} p(s', r | s, a) [r + \gamma V_*(s')]$$

$$\begin{aligned} q_*(s, a) &= E[R_{t+1} + \gamma V_{t+1} | S_t = s, A_t = a] \\ &= E[R_{t+1} + \gamma \max_{a'} q_*(s_{t+1}, a') | S_t = s, A_t = a] \end{aligned}$$

$$q_*(s, a) = \sum_{s', r} p(s', r | s, a) [r + \max_{a'} q_*(s', a')]$$

Q1.1 What is optimality?

→ Refers to best possible policy agent can follow to maximize expected return.

Q1.2 What is expectation & how does it relate to bellman eq<sup>n</sup>?

→ Expectation:- Mean of possible values of random Variable

$$E(x) = \sum p(x=x) \cdot x$$

Bellman eq<sup>n</sup> is recursive relation bet<sup>n</sup> value function which is expected return in a state

Q1.3 Solve exercise 3.25 to 3.29.

→ See sol<sup>n</sup> of Chapter 3 pdf in repo provided sol<sup>n</sup> in



# 1. Intuition behind Policy Iteration & Value Iteration

→ How are they working:-

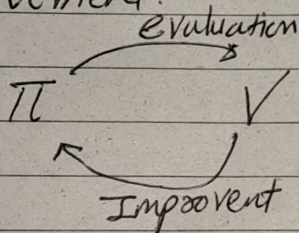
at every point you try out policy & evaluate it, based on evaluation policy is tweaked.

Why should it work

Ans → At every step either policy is kept same or improved. Hence we converge to optimal policy.

## 2. What is BPI?

→ Interaction of policy evaluation & policy improvement.



Loop continued till policy is stable.

Q4

Policy Iteration Time Complexity  $O(S^2A)$   
 Space  $O(S)$   
 Similar for Value Iteration

But Policy Iteration → Few Iterations  
 Value Iteration → More Iterations

For Small MDP → Policy Iteration  
 Large MDP → Value Iteration