Q1>



**MDP diagram**

| S | a | s' | $P(s', r \mid s, a)$ | $r(s, a, s')$ |
|---|---|---|---|---|
| Canteen | study | Canteen | 0.1 | +1 |
| Canteen | study | hostel | 0.3 | -1 |
| Canteen | study | Academic Build. | 0.6 | +3 |
| Canteen | eat | canteen | 1 | +1 |
| hostel | study | hostel | 0.5 | -1 |
| hostel | study | Academic Build. | 0.5 | +3 |
| hostel | eat | canteen | 1 | +1 |
| Academic Build | study | Academic Build | 0.7 | +3 |
| Academic Build | study | canteen | 0.3 | +1 |
| Academic Build | eat | Academic Build | 0.2 | +3 |
| Academic Build | eat | canteen | 0.8 | +1 |

other transitions have Probability = 0 and hence no reward.

⇒ Both value and policy iteration show same trend in state-value function,

⟹ V(academic building) > V(canteen) > V(hostel)

as expected because of the rewards associated with these places show same pattern.

⇒ The optimal policy in both cases was to 'study' in any locations when possible.