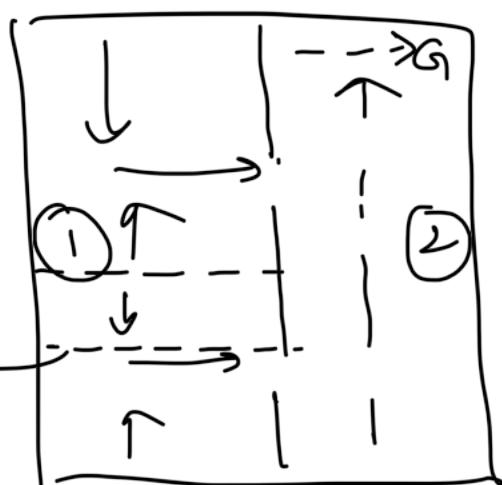


Wk - II RL

L1 Hierarchical RL

- provides temporal abstraction
- transfer / reusability of the funcⁿ
- more powerful / "meaningful" state abstraction

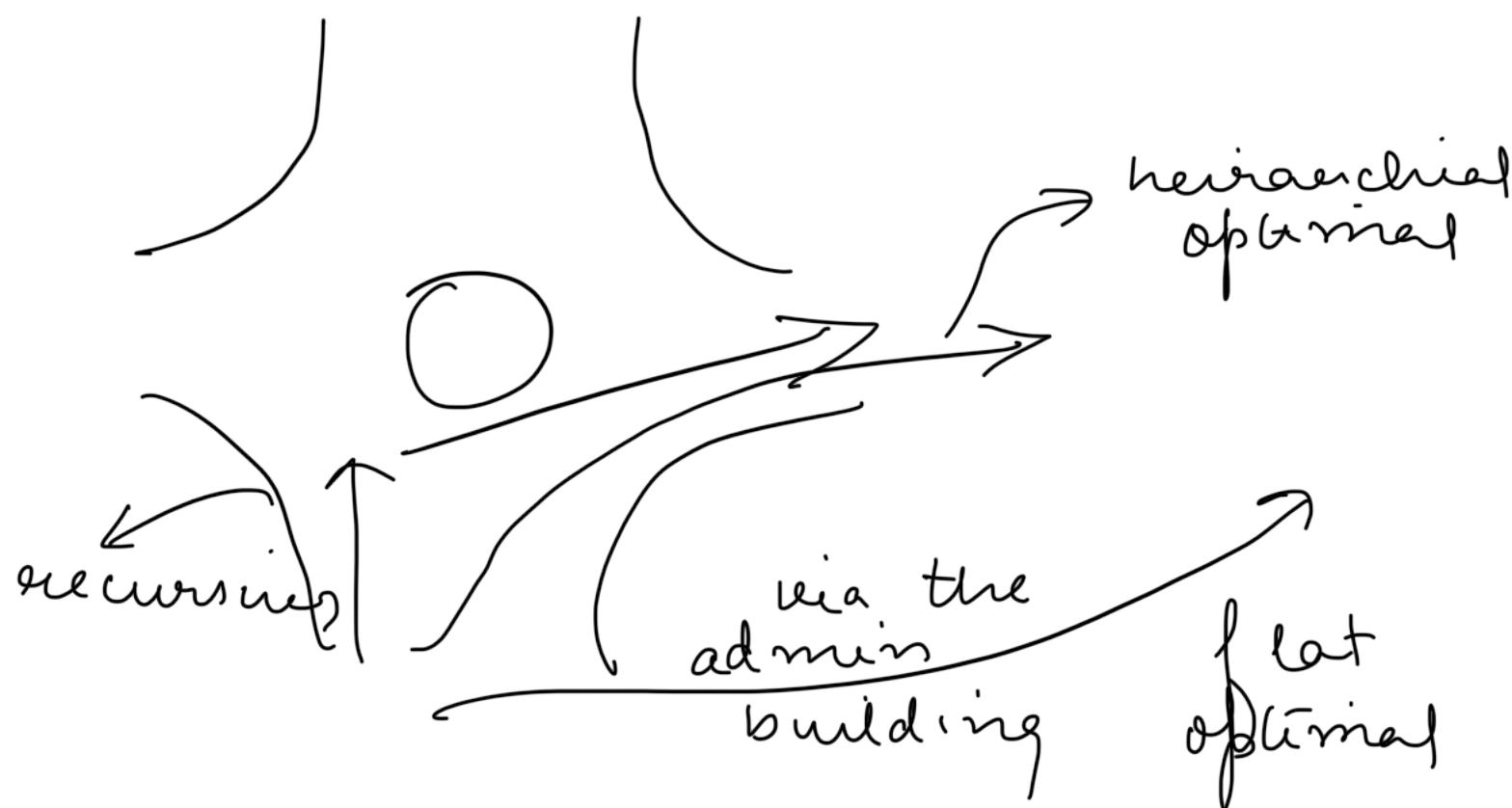
L2 Types of Optimality



then
the policy
will change

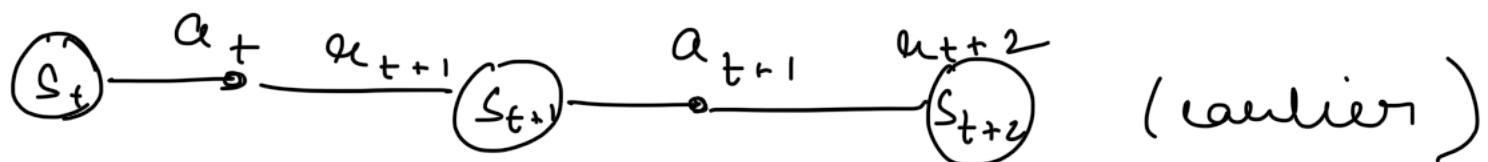
- there are many types of optimality -
 1. hierarchical optimal
 2. necessarily optimal (at least each subproblem is optimal)
 3. flat optimal

Eg. going to hostel zone via GIC.

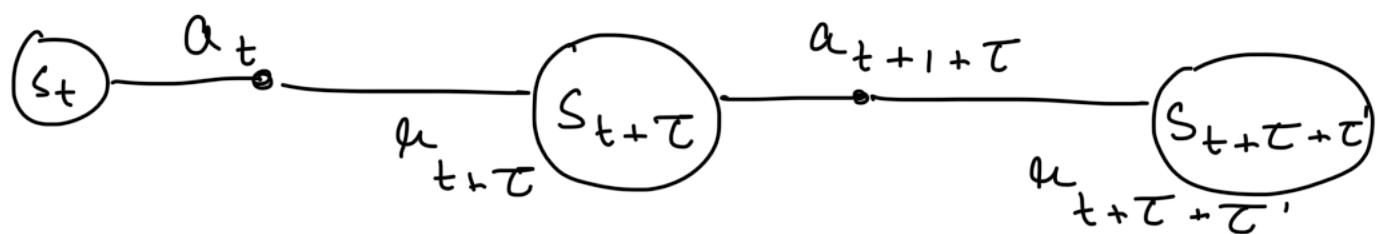


L3 Semi-Markov Decision Process

- SMDP (actions with durations)



now,



here, τ is referred to as the holding time

$$\langle S, A, P, R \rangle \quad p(s', \tau | s, a)$$

$$R = E[r | s, a, s', \tau]$$

- SMDP Q Learning

$$Q(s, a) = Q(s, a) + \alpha [r_{t+\tau} + \gamma^{t+\tau-1} \max_{a'} Q(s', a') - Q(s, a)]$$

$$u_{t+\tau} = u_{t+1} + \gamma u_{t+2} + \dots + \gamma^{\tau-1} u_{t+\tau}$$

LY Options

- It helps to choose macro actions . kind of sub-problem to the main problem.
 - We will be using the encapsulated policy concept . In this , basically we bring the policy for each sub-problem under one umbrella .
 - Initial Set , $I_0 \subseteq S$
 - Policy , $\pi_0 : S \rightarrow A$
 - Terminal $\beta : S \rightarrow [0, 1]$
- ∴ define the option as

$$O = \langle I_0, \pi_0, \beta \rangle$$

- Markov options - π_t , depends only on the current state.
- Semi-markov options - π_t , depends on the history, exactly, since when the option had started.

L5 Learning With Options

- SMDP Q learning
- $Q(s, a) \Rightarrow$ that we defined above
- Now, we come to Inter-options Q learning

$$Q(s, o) \quad s_1, \dots, s_n \\ \pi_o \rightarrow a_1, \dots, a_n$$

$$Q(s, a_1) = Q(s_1, a_1) + \alpha [r_1 + \gamma \max_a Q(s_1, a) \\ \text{similarly for,} \\ Q(s_1, a_2) \\ \vdots \\ - Q(s_1, a_1)]$$

$$Q(s_1, o) = Q(s_1, a) + \alpha [r_1 + \gamma Q(s_2, o) \\ - Q(s_1, o)] \\ \rightarrow \text{if not ending at } s_2.$$

if it terminates at s_2 , then

$$\gamma Q(s_2, 0) \Rightarrow \gamma \max_a Q(s_2, a)$$