

Lecture-11

Hierarchical RL

flat RL

The diagram illustrates the relationship between primitive and macro actions in a grid-based environment.

Primitive actions → N,S,E,W

Macro actions → exit the room from east door

The grid shows a room with a character at the bottom center. The room has doors on the top, right, and left sides. An arrow points upwards from the grid, indicating the direction of primitive actions (Up).

The text "exit the room from east door" is written above the grid, with arrows pointing to the right (East) and down (South). Below this, another set of arrows points left (West) and up (North), indicating the directions of macro actions.

Hierarchical RL

aims to find good re-usable temporally extended actions that may also provide opportunities for state abstractions.

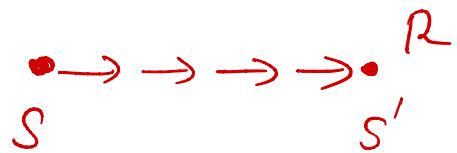
Advantages:

- ① Scalability
 - ② reusability / transfer learning,
 - ③ parallelize learning,
 - ④ long-term credit assignment
 - ⑤ structured exploration.

Abstract actions:-



- temporally extended actions



Ex: Room-leaving action.

- 'Macros' in programming

higher level RL agent \leftrightarrow Program.

Options/macros-action \leftrightarrow Sub-routines.

Semi-MDP :-

SMDP \rightarrow MDP that includes abstract actions -

$N \geq 1$ — random variable that denotes the # of steps that an abstract action 'a' takes to complete.

Transition fn: $T: S \times A \times S \times N \rightarrow [0, 1]$

$$T(s, a, s', N) = \Pr\{S_{t+N} = s' \mid S_t = s, a_t = a\}$$

Reward fn: $R: S \times A \times S \times N \rightarrow \mathbb{R}$

$$R(s, a, s', N) = \mathbb{E}\left[\sum_{n=0}^{N-1} \gamma^n r_{t+n} \mid S_t = s, a_t = a, S_{t+N} = s'\right]$$

State value fn

$$V^\pi(s) = \mathbb{E}_\pi \left[\sum_{k=0}^{\infty} \gamma^k r_{t+k} \mid s_t = s \right]$$

$$= \mathbb{E}_\pi \left[(r_t + \gamma r_{t+1} + \dots + \gamma^{N-1} r_{t+N-1}) + \right. \\ \left. (\gamma^N r_{t+N} + \dots) \right] \text{Step}$$

$$V_\pi(s) = \sum_{s', N} T(s, \pi(s), s', N) \left[R(s, \pi(s), s', N) + \gamma^N V^\pi(s') \right]$$

Optimal value fn

$$v^*(s) = \max_a \sum_{s', N} T(s, a, s', N) \left[R(s, a, s', N) + \gamma^N v^*(s') \right]$$

State action value fn

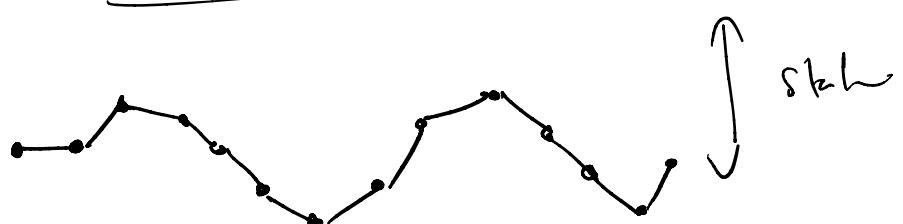
$$Q^\pi(s, a) = \sum_{s', N} T(s, a, s', N) \left[R(s, a, s', N) + \gamma^N Q^\pi(s', \pi(s')) \right]$$

Optimal Q-fn

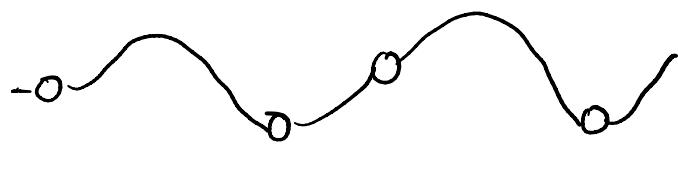
$$Q^*(s, a) = \sum_{s', N} T(s, a, s', N) \left[R(s, a, s', N) + \gamma^N \max_{a'} Q^*(s', a') \right]$$

time →

MDPs



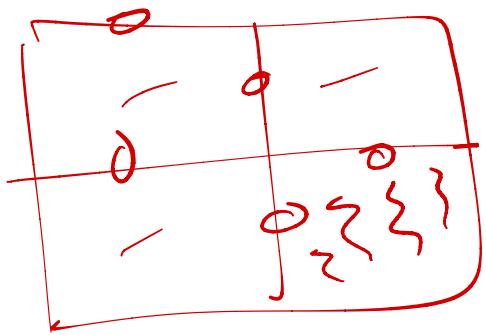
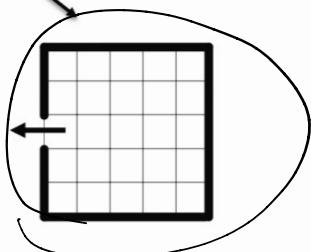
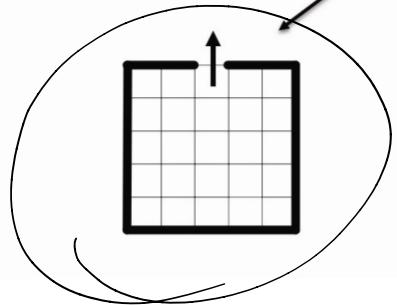
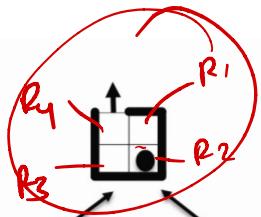
SMDPs



— Structure

North room-leaving abstract action

West room-leaving abstract action

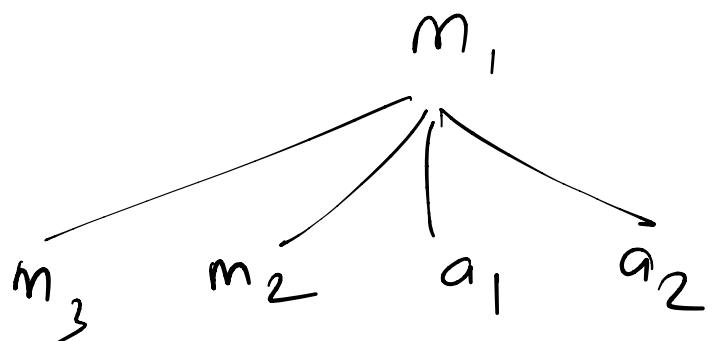


Optimality: \rightarrow flat optimality ✓

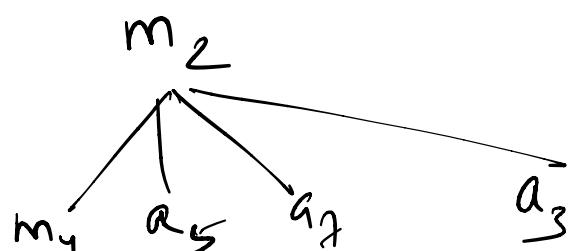
\rightarrow hierarchical optimality ✓

\searrow Recursive optimality

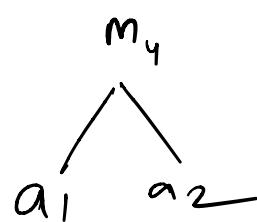
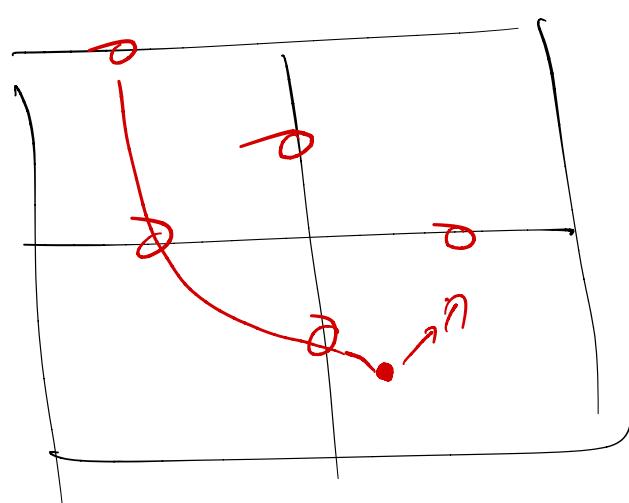
$m_1 : a_1, a_2$



$m_1 : m_2, a_1, a_2$

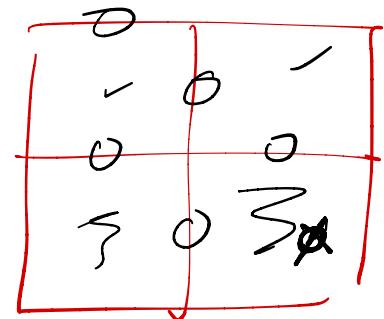


$m_1 : a_5, a_1, a_2$



① Options (Sutton et al. 1999) :

MDP: $\langle S, A, T, R \rangle$



Option: $\langle I, \pi, \beta \rangle$

$I \subseteq S$ initiation set.

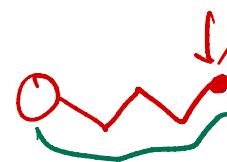
$\pi : S \times A \rightarrow [0, 1]$ Policy.

$\beta : S^+ \rightarrow [0, 1]$ is a termination condn.

Markov options



Semi-Markov options



Primitive actions are also options.

a : $\langle I = s, \pi(s, a) = 1.0, \beta(s) = 1 \rangle$

for all $s \in S$.

SMDP Q-learning :-

Cum.
reward by
end of the
option

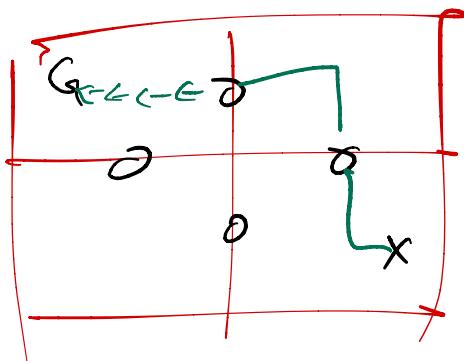
$$Q(s, o) \leftarrow Q(s, o) + \alpha \left[r + \gamma^k \max_{a \in O} Q(s', a) \right]$$

$\rightarrow Q(s, o)$

O - option

$k \rightarrow$ # of timesteps between s, s' .

r - Cumulative disc. reward over this time pen —



$O \Rightarrow \{$

- $N, S, E, W,$
- exit thru N
- exp. in S
- exp. in E
- exp. in W

Intra-option Q learning :-

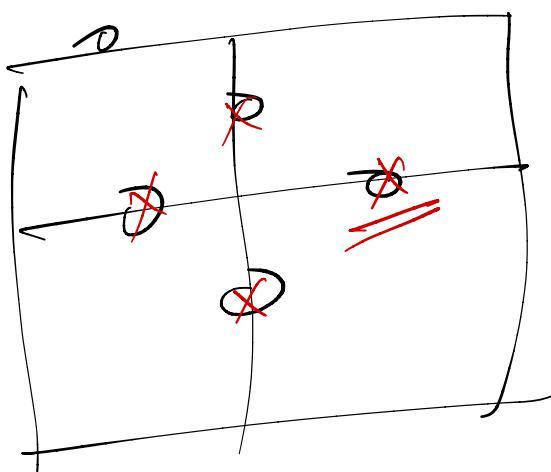
$$\tilde{Q}_o^*(s, o) = (1 - \beta(s)) Q_o^*(s, o) + \beta(s) \max_{o' \in O} Q_o^*(s, o')$$

$$Q_o^*(s, o) = \sum_{a \in A_s} \pi(s, a) E \left[r + \gamma \tilde{Q}_o^*(s', o') \mid s, a \right]$$

One-step TD update :-

$$Q(s_t, o) \leftarrow Q(s_t, o) + \alpha$$

$$\left[r_{t+1} + \gamma \tilde{Q}_o^*(s_{t+1}, o) - Q(s_t, o) \right]$$

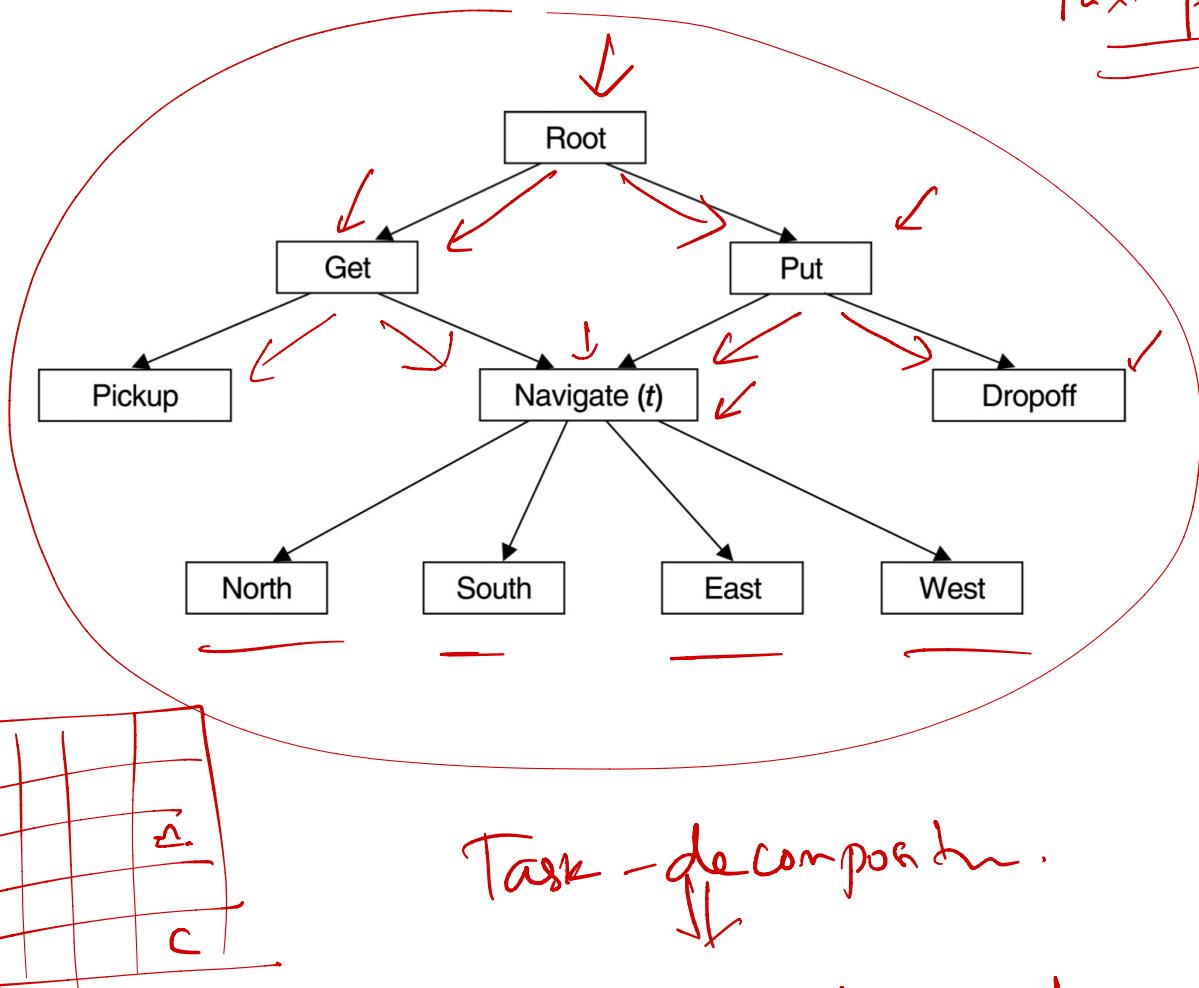


"bottle neck state".

MAX Q

(Dietterich, 2000)

Taxi problem



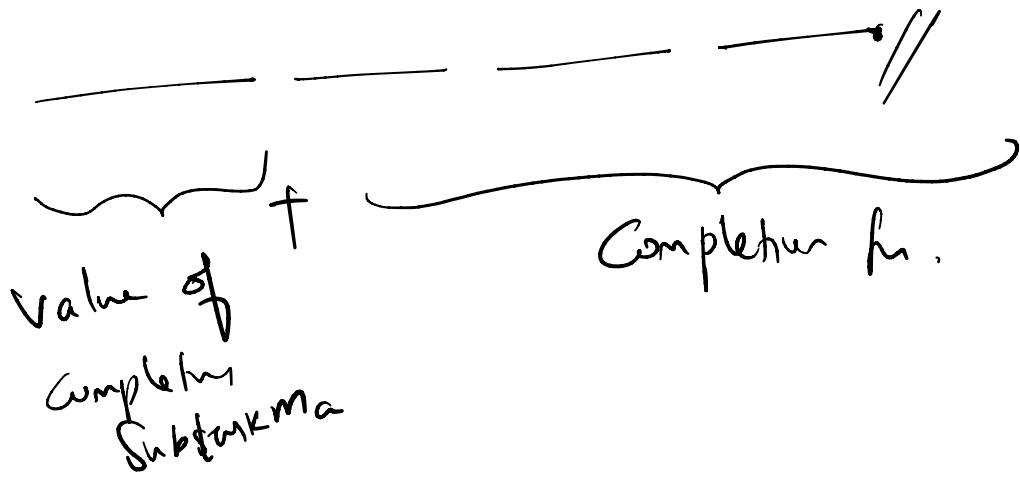
Task-decomposition.

value for decomp.

m-Subtask

$$Q^{\pi}(m, s, a) = \sum_{s', N} T(m, s, a, s', N)$$

$$\left[R(m, s, a, s', N) + \gamma^N Q^{\pi}(m, s', \pi(s')) \right]$$

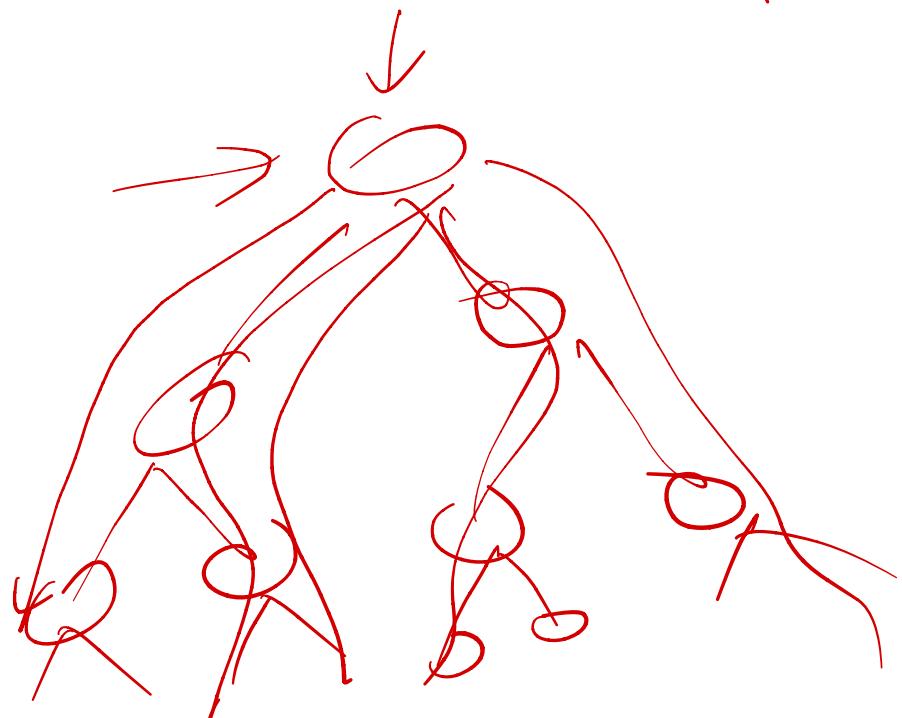


$$C^{\pi}(m, s, a) = \sum_{s', N} T(m, s, a, s', N) \gamma^N Q^{\pi}(m, s', \pi(s'))$$

$$Q^{\pi}(m, s, a) = V^{\pi}(m_a, s) + C^{\pi}(m, s, a)$$

$$V^{\pi}(m_a, s) = \begin{cases} \sum_{s'} T(s, a, s') R(s, a, s') & \text{if } a \text{ is pm.} \\ Q^{\pi}(m_a, s, \pi(s)) & \text{if } a \text{ is abfl.} \end{cases}$$

$$Q^{\pi}(m_0, s, \pi(s)) = r^{\pi}(m_k, s) + C^{\pi}(m_{k-1}, s, a_{k-1}) \\ \dots + C^{\pi}(m_1, s, a_1) + C^{\pi}(m_0, s, a_0)$$



$\max Q$,

$$r^{\pi}(m_a, s) = \max_{a'} Q^{\pi}(m_a, s, a')$$

Require: V, C, A , learning rate α

```

1: if  $m$  is a primitive action then
2:   execute  $m$ , receive reward  $r$ , and observe the result state  $s'$ 
3:    $V(m,s) \leftarrow (1 - \alpha)V(m,s) + \alpha r$ 
4:   return 1
5: else
6:    $count = 0$ 
7:   while  $m$  has not terminated do
8:     choose an exploration action  $a$ 
9:      $N = MAXQ(a,s)$ 
10:    observe result state  $s'$ 
11:     $\langle V(j^{greedy}, s'), a_{j^{greedy}} \rangle \leftarrow Evaluate(m, s')$ 
12:     $C(m,s,a) \leftarrow (1 - \alpha)C(m,s,a) + \alpha \gamma^N V(j^{greedy}, s')$ 
13:     $count \leftarrow count + N$ 
14:     $s = s'$ 
15:   return  $count$ 

```

Algorithm 19. $MAXQ(m,s)$ [Dietterich (2000)]

