

MODULAR MULTITASK RL w/ POLICY SKETCHES

If we give RL agents a description of high-level actions to take in order to compute a complex task, w/o specifying the low-level actions that comprise the high level actions, can the agent learn those low-level associations?

GIVEN H_b , the set of high level tasks.

Provide agent X (which) with list of tasks $h \in H_b$, can it learn low level actions & policy associated with every h ?

Structured Deep Models ← Compositional reasoning & interaction

Used for QA systems
Relational reasoning

Page 3 P 3 → not sure about the task-specific training signals.

NCP

Natural lang.

instruction-following: string → symbolic action → Hand-coded sequences → relations

THIS PAPER

Symbiosis Action Sequences → Learned Action Policies [discovering on own]

3. LEARNING MODULAR POLICIES FROM SKETCHES.

(S, A, P, γ)

S : set of states

A : set of LOWLEVEL actions.

$P: S \times A \times S \rightarrow R$
transition prob. dist.

γ : discount factor

Standard inf. horizon
discounted MDP

task $T \in T$ specified by (R_T, π_T) :

$R_T: S \rightarrow R$ (task specific reward function)

$\pi_T: S \rightarrow R$ initial distribution over states.

$\{(s_i, a_i)\}$ sequence of states & actions
 $q_i = \sum_{j=i+1}^{\infty} \gamma^{j-i-1} R(s_j)$ empirical return.

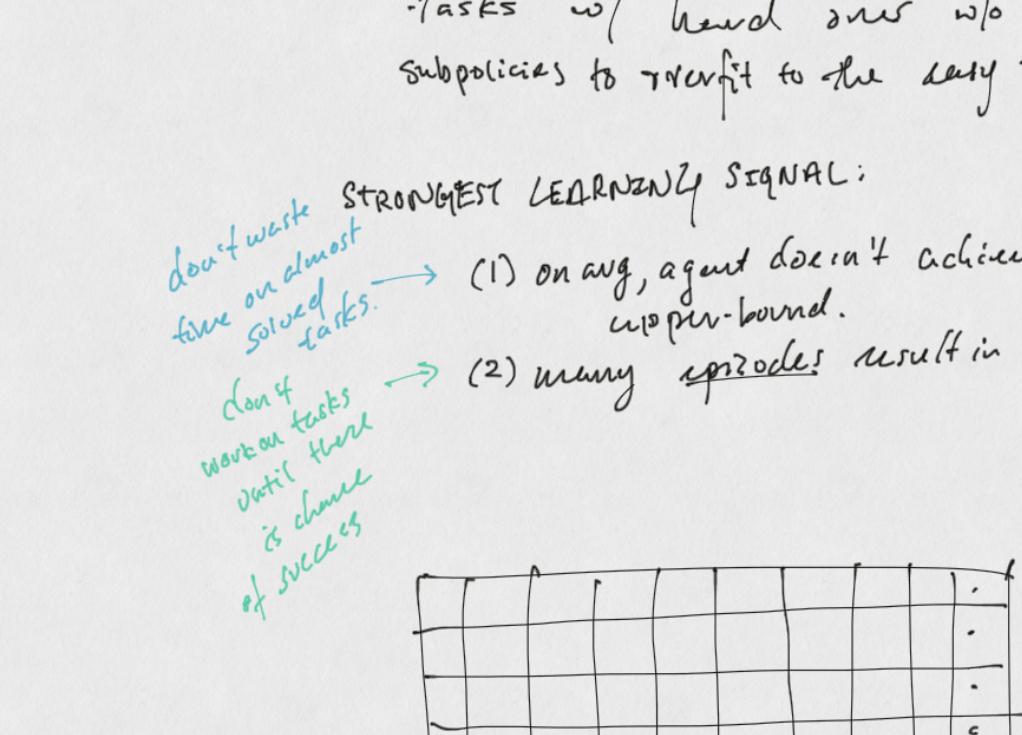
ANNOOTE TASKS $T \in T$ with sketches K_T :

$K_T = (b_{T1}, b_{T2}, \dots)$
 $b_T \in \mathcal{B}$ fixed vocab of symbolic labels

3.1 MODEL

construct for each symbol b , subpolicy π_b .

$$T_1 = (b_{T1}, b_{T2}, \dots) \quad T_2 = (b_{T2}, b_{T2}, \dots)$$



At each timestep,

π_b can choose A or STOP

$$\Rightarrow \mathcal{A}^+ = \mathcal{A} \cup \{\text{STOP}\}$$

$\pi_b: S \rightarrow \mathcal{A}^+$ ANY REPRESENTATION WORKS.
 $\rightsquigarrow \pi$ is NN in THIS PAPER

$$K_T = (b_1, b_2, \dots) \rightarrow \pi_T = \text{CONCAT}(\pi_{b_1}, \pi_{b_2}, \dots)$$

Maintain subpolicy index i
when π_{b_i} emits STOP , inc. i .

$\pi = \bigcup_T \{\pi_T\}$, π_T is arbitrary func. $S \rightarrow \mathcal{A}$ param. by θ_b .

optimize over all θ_b

$$J(\pi) = \sum_T J(\pi_T) = \sum_T \mathbb{E}_{s_i \sim \pi_T} \left[\sum_i \gamma^i R_T(s_i) \right]$$

maximize expected discounted reward

3.2 POLICY OPTIMIZATION

$$\nabla_{\theta} J(\pi) = \sum_i (\nabla_{\theta} \log \pi(a_i | s_i)) (q_i - c(s_i))$$

one critic per TASK T

subpolicies partake in multiple tasks w/ different reward functions

\Rightarrow no well defined subpolicy STATE-VALUE function.

3.3 CURRICULUM LEARNING

Complex tasks → impossible to earn ANY positive reward w/o already establishing some subpolicies.

Curriculum learning balances "EASY"

Tasks w/ hard ones w/o allowing subpolicies to transit to the easy tasks.

don't waste STRONGEST LEARNING SIGNAL:
take on easiest solved tasks. → (1) on avg, agent doesn't achieve new upper-bound.

don't want tasks until hard is done
is done
is done

many episodes result in high reward

do?

$E[H_{12}] = \frac{1}{2} E[H_{11}] + \frac{1}{2} E[H_{13}]$

$$H_{12} = H_{13}$$

$$H_{13} = \frac{1}{2} H_{12} + \frac{1}{2} H_{17}$$

$$H_{17} = \frac{1}{2} H_{18} + \frac{1}{2} H_{16}$$

$$H_{18} = \frac{1}{2} H_{19} + \frac{1}{2} H_2$$

$$H_2 = \frac{1}{2} H_3 + \frac{1}{2} H_1$$

$$H_3 = H_2$$

$$H_2 = \frac{1}{2} (H_3 + 1) + \frac{1}{2}$$

$$\frac{1}{2} H_2 = 1$$

$$H_2 = 2$$

$$H_2 = \frac{1}{2} (H_3 + 2)$$

$$H_3 = \frac{1}{2} H_2 + \frac{1}{2} + 1 \Rightarrow H_3 = \frac{1}{2} (\frac{1}{2} H_2 + 2)$$

$$H_2 = \frac{1}{2} H_3 + \frac{1}{2} + 2 = \frac{1}{2} H_3 + \frac{2}{3} (\frac{1}{2} H_2 + 2) + 1$$

$$H_3 = 1 + H_2 \Rightarrow H_3 = \frac{3}{2} (\frac{1}{2} H_2 + \frac{2}{3})$$

$$H_2 = 1 + \frac{3}{2} (\frac{1}{2} H_2 + \frac{2}{3})$$

$$\frac{1}{2} H_2 = 1 + \frac{7}{2} = \frac{9}{2}$$

$$H_2 = 18$$

$$H_1 = 18$$

$$H_1 = 1 + \frac{3}{2} (\frac{1}{2} H_2 + \frac{2}{3})$$

$$\frac{1}{2} H_1 = 1 + \frac{7}{2} = \frac{9}{2}$$

$$H_1 = 18$$

$$H_1 = 1 + \frac{3}{2} (\frac{1}{2} H_2 + \frac{2}{3})$$

$$\frac{1}{2} H_1 = 1 + \frac{7}{2} = \frac{9}{2}$$

$$H_1 = 18$$

$$H_1 = 1 + \frac{3}{2} (\frac{1}{2} H_2 + \frac{2}{3})$$

$$\frac{1}{2} H_1 = 1 + \frac{7}{2} = \frac{9}{2}$$

$$H_1 = 18$$

$$H_1 = 1 + \frac{3}{2} (\frac{1}{2} H_2 + \frac{2}{3})$$

$$\frac{1}{2} H_1 = 1 + \frac{7}{2} = \frac{9}{2}$$

$$H_1 = 18$$

$$H_1 = 1 + \frac{3}{2} (\frac{1}{2} H_2 + \frac{2}{3})$$

$$\frac{1}{2} H_1 = 1 + \frac{7}{2} = \frac{9}{2}$$

$$H_1 = 18$$

$$H_1 = 1 + \frac{3}{2} (\frac{1}{2} H_2 + \frac{2}{3})$$

$$\frac{1}{2} H_1 = 1 + \frac{7}{2} = \frac{9}{2}$$

$$H_1 = 18$$

$$H_1 = 1 + \frac{3}{2} (\frac{1}{2} H_2 + \frac{2}{3})$$

$$\frac{1}{2} H_1 = 1 + \frac{7}{2} = \frac{9}{2}$$

$$H_1 = 18$$

$$H_1 = 1 + \frac{3}{2} (\frac{1}{2} H_2 + \frac{2}{3})$$

$$\frac{1}{2} H_1 = 1 + \frac{7}{2} = \frac{9}{2}$$

$$H_1 = 18$$

$$H_1 = 1 + \frac{3}{2} (\frac{1}{2} H_2 + \frac{2}{3})$$

$$\frac{1}{2} H_1 = 1 + \frac{7}{2} = \frac{9}{2}$$

$$H_1 = 18$$

$$H_1 = 1 + \frac{3}{2} (\frac{1}{2} H_2 + \frac{2}{3})$$

$$\frac{1}{2} H_1 = 1 + \frac{7}{2} = \frac{9}{2}$$

$$H_1 = 18$$

$$H_1 = 1 + \frac{3}{2} (\frac{1}{2} H_2 + \frac{2}{3})$$

$$\frac{1}{2} H_1 = 1 + \frac{7}{2} = \frac{9}{2}$$

$$H_1 = 18$$

$$H_1 = 1 + \frac{3}{2} (\frac{1}{2} H_2 + \frac{2}{3})$$

$$\frac{1}{2} H_1 = 1 + \frac{7}{2} = \frac{9}{2}$$

$$H_1 = 18$$

$$H_1 = 1 + \frac{3}{2} (\frac{1}{2} H_2 + \frac{2}{3})$$

$$\frac{1}{2} H_1 = 1 + \frac{7}{2} = \frac{9}{2}$$

$$H_1 = 18$$

$$H_1 = 1 + \frac{3}{2} (\frac{1}{2} H_2 + \frac{2}{3})$$

$$\frac{1}{2} H_1 = 1 + \frac{7}{2} = \frac{9}{2}$$

$$H_1 = 18$$

$$H_1 = 1 + \frac{3}{2} (\frac{1}{2} H_2 + \frac{2}{3})$$

$$\frac{1}{2} H_1 = 1 + \frac{7}{2} = \frac{9}{2}$$

$$H_1 = 18$$

$$H_1 = 1 + \frac{3}{2} (\frac{1}{2} H_2 + \frac{2}{3})$$

$$\frac{1}{2} H_1 = 1 + \frac{7}{2} = \frac{9}{2}$$

$$H_1 = 18$$

$$H_1 = 1 + \frac{3}{2} (\frac{1}{2} H_2 + \frac{2}{3})$$

$$\frac{1}{2} H_1 = 1 + \frac{7}{2} = \frac{9}{2}$$

$$H_1 = 18$$

$$H_1 = 1 + \frac{3}{2} (\frac{1}{2} H_2 + \frac{2}{3})$$

$$\frac{1}{2} H_1 = 1 + \frac{7}{2} = \frac{9}{2}$$

$$H_1 = 18$$

$$H_1 = 1 + \frac{3}{2} (\frac{1}{2} H_2 + \frac{2}{3})$$

$$\frac{1}{2} H_1 = 1 + \frac{7}{2} = \frac{9}{2}$$

$$H_1 = 18$$