



OPTION DISCOVERY USING DEEP SKILL CHAINING

Akhil Bagaria

Department of Computer Science Brown University Providence, RI, USA akhil_bagaria@brown.edu George Konidaris Department of Computer Science Brown University Providence, RI, USA gdk@brown.edu

Presenter: LinusWangg

Background



Problems:

1. Simple algorithms suffer from the pressure of the feature representation.

2. Long-horizon tasks is hard to handle.

3. Sparse Reward.



Alpha Go



FPS game

Literature reviews



Algorithms and Architecture



Option-Critic Architecture



HAC

Method-Deep Skill Chaining



Implements:

- 1. Based on the Option-Critic Architecture.
- 2. Using the Neural Network.
- 3. Extension of the original Skill Chaining.



Implement Graph

Method-Deep Skill Chaining



Graph illustration:



Method-Deep Skill Chaining



Option Select:

$$\mathcal{O}'(s_t) = \{ o_i | \mathcal{I}_{o_i}(s_t) = 1 \cap \beta_{o_i}(s_t) = 0, \forall o_i \in \mathcal{O} \}$$

$$\tag{2}$$

$$o_t = \underset{o_i \in \mathcal{O}'(s_t)}{\arg \max} Q_{\phi}(s_t, o_i).$$
(3)

Option Value Update:





Deep Q Learning

Experiment

ParNeC 模式识别与神经计算研究组 PAttern Recognition and NEural Computing



Initialization Set Hot Map



Figure 3: Solution trajectories found by deep skill chaining. Sub-figure (d) shows two trajectories corresponding to the two possible initial locations in this task. Black points denote states in which $\pi_{\mathcal{O}}$ chose primitive actions, other colors denote temporally extended option executions.

Individual Option Performance

Experiment





Reward Results

Discussions



Drawbacks:

1. This method require the agent to reach the goal state without any prior knowledge, it's too hard!

2. the gestation period of option is affected by the exploration.



 $s_0, a_0, s_1, a_1, \ldots, s_t, a_t, \ldots, s_{t+H-1}, a_{t+H-1}, s_{t+H}, a_{t+H}, s_{t+H+1}, a_{t+H+1}, \ldots, s_T$

Discussions

Solutions:

1. Skill extractions from the existing trajectories in offline RL or IL.

2. Use some more efficient explore policies.



(a) Previous skill-based imitation (b) Previous skill-based imitation (c) Our method with subtask-local with demonstration-covering dataset with subtask-local dataset

Imitation Learning

1. Intrinsic reward (green) is distributed throughout the environment







2. An IM algorithm might start by exploring (purple) a nearby area with intrinsic reward

模式识别与神经计算研究组

PAttern Recognition and NEural Computing



4. Exploration fails to rediscover promising areas it has detached from



Go explore

ParN_pC





Architecture:

Option-Critic Architecture

Problems solved:

- 1. Automatically separate the long-horizon task into options.
- 2. Relieve the pressure of feature representation.
- 3. Realize the generalization of the Skill Chaining building.

Potential improvement:

- 1. Combining with some offline RL and IL algorithms.
- 2. Some good explore policies.

Thanks