

Modular Multitask Reinforcement Learning with Policy Sketches

Yoonho Lee

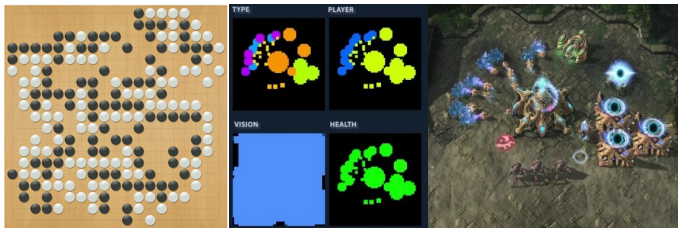
Department of Computer Science and Engineering
Pohang University of Science and Technology

August 09, 2017

Modular Multitask Reinforcement Learning with Policy Sketches

- ▶ ICML 2017 Best Paper Honorable Mention
- ▶ Hierarchical Reinforcement Learning
- ▶ Proposes a looser form of task supervision for RL agents (compared to e.g. reward shaping)
- ▶ Proposed form of task supervision is decoupled from the environment
- ▶ Proposes a new way to use NNs for hierarchical RL
- ▶ Natural extensions to zero-shot and unlabelled RL

Hierarchical Reinforcement Learning



- ▶ Why is starcraft harder than go for RL?
- ▶ Why is starcraft not harder than go for humans?

Hierarchical Reinforcement Learning

- ▶ Decision \neq Action. RL algorithms are designed with decisions in mind, but operate on actions.
- ▶ Real world decisions have hierarchical structure(e.g. cook dinner \rightarrow cut potato \rightarrow activate arm muscle)
- ▶ Effective knowledge reuse
- ▶ Efficient credit assignment
- ▶ Open question: How do we discover salient/reusable decisions?

Previous Work

- ▶ Learn multiple timestep policy along with confidence¹
- ▶ Learn controller network that sets desired direction of state change²
- ▶ Encode multiple sub-policies with an SNN³
- ▶ Actor-Critic algorithm for options⁴

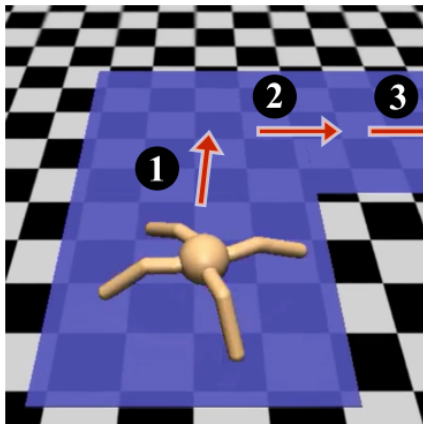
¹Alexander Vezhnevets et al. “Strategic Attentive Writer for Learning Macro-Actions”. In: *NIPS* (2016).

²Alexander Sasha Vezhnevets et al. “FeUdal Networks for Hierarchical Reinforcement Learning”. In: (2017).

³Carlos Florensa, Yan Duan, and Pieter Abbeel. “Stochastic Neural Networks for Hierarchical Reinforcement Learning”. In: *ICLR 2017* (2017), pp. 1056–1064.

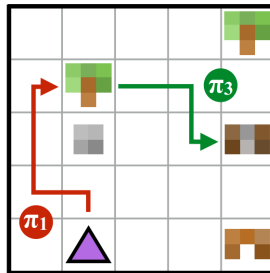
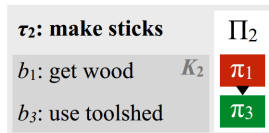
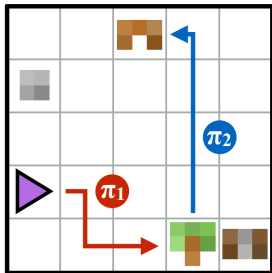
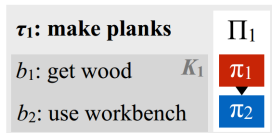
⁴Pierre-Luc Bacon, Jean Harb, and Doina Precup. “The Option-Critic Architecture”. In: *AAAI* (2017).

Previous Work



- Previous hierarchical RL algorithms require reward engineering for high-dimensional environments

Proposed Approach

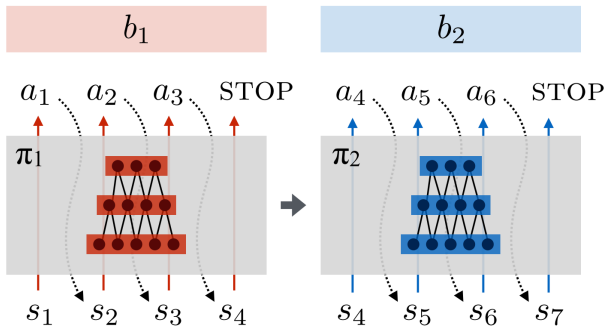


Proposed Approach

Sketches

Goal	Sketch				
Crafting environment					
make plank	get wood	use toolshed			
make stick	get wood	use workbench			
make cloth	get grass	use factory			
make rope	get grass	use toolshed			
make bridge	get iron	get wood	use factory		
make bed*	get wood	use toolshed	get grass	use workbench	
make axe*	get wood	use workbench	get iron	use toolshed	
make shears	get wood	use workbench	get iron	use workbench	
get gold	get iron	get wood	use factory	use bridge	
get gem	get wood	use workbench	get iron	use toolshed	use axe

Method



Method

Algorithm 1 TRAIN-STEP(Π , curriculum)

```
1:  $\mathcal{D} \leftarrow \emptyset$ 
2: while  $|\mathcal{D}| < D$  do
3:   // sample task  $\tau$  from curriculum (Section 3.3)
4:    $\tau \sim \text{curriculum}(\cdot)$ 
5:   // do rollout
6:    $d = \{(s_i, a_i, (b_i = K_{\tau,i}), q_i, \tau), \dots\} \sim \Pi_\tau$ 
7:    $\mathcal{D} \leftarrow \mathcal{D} \cup d$ 
8:   // update parameters
9:   for  $b \in \mathcal{B}, \tau \in \mathcal{T}$  do
10:     $d = \{(s_i, a_i, b', q_i, \tau') \in \mathcal{D} : b' = b, \tau' = \tau\}$ 
11:    // update subpolicy
12:     $\theta_b \leftarrow \theta_b + \frac{\alpha}{D} \sum_d (\nabla \log \pi_b(a_i | s_i)) (q_i - c_\tau(s_i))$ 
13:    // update critic
14:     $\eta_\tau \leftarrow \eta_\tau + \frac{\beta}{D} \sum_d (\nabla c_\tau(s_i)) (q_i - c_\tau(s_i))$ 
```

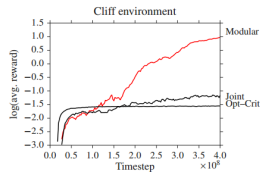
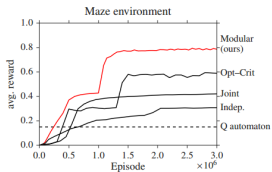
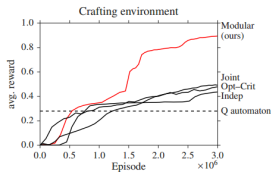
Method

Algorithm 2 TRAIN-LOOP()

```
1: // initialize subpolicies randomly
2:  $\Pi = \text{INIT}()$ 
3:  $\ell_{\max} \leftarrow 1$ 
4: loop
5:    $r_{\min} \leftarrow -\infty$ 
6:   // initialize  $\ell_{\max}$ -step curriculum uniformly
7:    $\mathcal{T}' = \{\tau \in \mathcal{T} : |K_{\tau}| \leq \ell_{\max}\}$ 
8:    $\text{curriculum}(\cdot) = \text{Unif}(\mathcal{T}')$ 
9:   while  $r_{\min} < r_{\text{good}}$  do
10:    // update parameters (Algorithm 1)
11:     $\text{TRAIN-STEP}(\Pi, \text{curriculum})$ 
12:     $\text{curriculum}(\tau) \propto \mathbb{1}[\tau \in \mathcal{T}'](1 - \hat{\mathbb{E}}r_{\tau}) \quad \forall \tau \in \mathcal{T}$ 
13:     $r_{\min} \leftarrow \min_{\tau \in \mathcal{T}'} \hat{\mathbb{E}}r_{\tau}$ 
14:     $\ell_{\max} \leftarrow \ell_{\max} + 1$ 
```

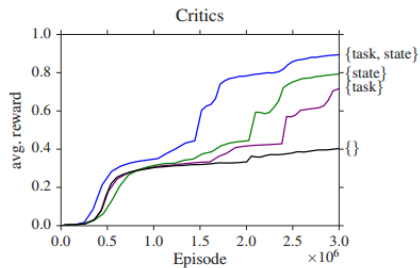
Experiments

Comparison to Baseline

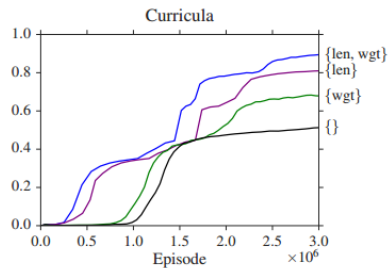


Experiments

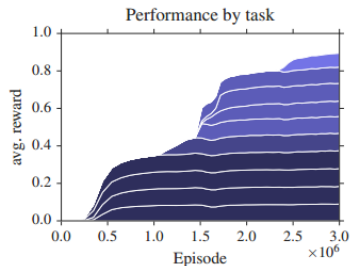
Ablation



(a)



(b)



Experiments

Accuracy

Model	Multitask	0-shot	Adaptation
Joint	.49	.01	–
Independent	.44	–	.01
Option–Critic	.47	–	.42
Modular (ours)	.89	.77	.76

Discussion

Weaknesses

- ▶ Only works for tasks that are strictly composed
- ▶ To label, one needs to understand which tasks are "RL-able"
- ▶ Labels needed: Fundamentally limited to tasks understood by humans

Discussion

Implications

- ▶ Zero-shot and Adaptation experiments had similar performance, what does this imply?
- ▶ Weight sharing between policies?
- ▶ Curriculum instead of policy sketches?
- ▶ Common theme of cooperating networks for hierarchical RL: Manager-Worker⁵, Option-Policy⁶, Policy modules⁷

⁵Alexander Sasha Vezhnevets et al. “FeUdal Networks for Hierarchical Reinforcement Learning”. In: (2017).

⁶Pierre-Luc Bacon, Jean Harb, and Doina Precup. “The Option-Critic Architecture”. In: AAAI (2017).

⁷Jacob Andreas, Dan Klein, and Sergey Levine. “Modular Multitask Reinforcement Learning with Policy Sketches”. In: ICML (2017). arXiv: 1611.01796.

References I

- [1] Jacob Andreas, Dan Klein, and Sergey Levine. “Modular Multitask Reinforcement Learning with Policy Sketches”. In: *ICML* (2017). arXiv: 1611.01796.
- [2] Pierre-Luc Bacon, Jean Harb, and Doina Precup. “The Option-Critic Architecture”. In: *AAAI* (2017).
- [3] Carlos Florensa, Yan Duan, and Pieter Abbeel. “Stochastic Neural Networks for Hierarchical Reinforcement Learning”. In: *ICLR 2017* (2017), pp. 1056–1064.
- [4] Alexander Vezhnevets et al. “Strategic Attentive Writer for Learning Macro-Actions”. In: *NIPS* (2016).
- [5] Alexander Sasha Vezhnevets et al. “FeUdal Networks for Hierarchical Reinforcement Learning”. In: (2017).

Thank You