
Attention-Over-Actions Option-Critic

Computer Science Extended Essay
Word Count:

Research Question

How are localized options trained in the Option-Critic architecture?

1 Introduction

As human we operate in high level actions. For example when driving a car, we make decisions about turning left or right instead of thinking about which muscle to contract. Human have the ability to group a chain of actions into one single high level action.

In Reinforcement Learning, we have a way to capture this idea of grouping action by using options [1]. When the options are defined, learning how to use them are very simple. However, if the options are not given and need to be learned, things get a lot harder since it requires knowing what makes an option good.

Many people argue that a good options should be diverse and localized, and many recent algorithms have followed this argument. Indeed, these algorithms have been able to discover options that are localized, so in this essay, I will try to derive a framework for localization in order to answer my research question, and I will develop an algorithm out of the framework to see whether it can actually train localized options.

This essay will be structured as follows: First, preliminary and related work will be presented to give a context to what I am trying to do. Second, previous work will be analyzed to figure out how localization is achieved. Third, a framework will be proposed based on the observations made in the analysis. Forth, an algorithm will be derived from the framework. Finally, the algorithm will be tested in the Four Rooms environment.

2 Preliminary

This section only acts as a summary. You are assumed to have basic knowledge about Reinforcement Learning and Options.

Markov Decision Process

Markov decision process (MDP) [4] is a mathematical framework for modeling decision making in a stochastic environment. It is defined as a tuple: $\langle \mathcal{S}, \mathcal{A}, r, \gamma, P \rangle$ where:

\mathcal{S} is the set of states.

\mathcal{A} is the set of actions

$r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is the reward function

$\gamma \in [0, 1)$ is the discount factor that ensure the cumulative reward $\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t)$ converges

$P : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ is the transition model which gives the probability for a particular transition to occur.

All MDPs must follow the Markov Property, which means that everything is stateless and does not depend on the history. In an MDP, A policy $\pi : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ is responsible for choosing the action, after an action is chosen, the environment transitions to a new state according to $P(s, a, s')$, and a reward is given to the policy. This process repeats until the environment enters a terminal state.

Reinforcement Learning

Reinforcement Learning (RL) [4] is a machine learning paradigm which allows an agent to learn from interaction in an MDP. The agent’s goal is to maximize a certain objective function.

Actor-Critic [5] is one of the popular classes of algorithms that harvest the advantage of both Q-Learning and Policy Gradient. An actor network is trained to choose the best action, while a critic network is trained to evaluate the decision made by the actor network.

Option Framework

In the Option Framework[1], instead of only using 1 policy, we use a set of options, each option is defined as a tuple: $\langle \mathcal{I}_\omega, \pi_\omega, \beta_\omega \rangle$, where:

$\mathcal{I}_\omega \subseteq \mathcal{S}$ is the initiation set that define which state the option can be selected

$\pi_\omega : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ is the internal policy

$\beta_\omega : \mathcal{S} \rightarrow [0, 1]$ is the termination probability.

These options take turns to choose the action. A policy-over-options $\pi_\Omega : \mathcal{S} \times \Omega \rightarrow [0, 1]$, where Ω is the set of options, decides which option to use. When an option is chosen, actions are chosen by its internal policy π_ω from then on, until it terminated according to its β_ω , then the policy-over-options π_Ω chooses an option again. An option has a chance to terminate every time environment transitions to a new state.

If the options are defined, the policy-over-options π_Ω can be learned by using SMDP Q-Learning [6] or Intra-Option Learning [7]. However, options are not always predefined and need to be discovered.

Four Rooms Environment

The Four Rooms environment is commonly used to evaluate algorithms that use options. The agent have 4 actions: up, down, left and right, which will move the agent in that corresponding direction. The choice of agent has 1/3 chance to fail and a random action will be chosen instead.

+50 reward will be given to the agent when it arrives to the goal state, and the episode will terminate after that.

Each cell position in the environment is mapped to a state number and is given to the agent in each step. The agent starts in a random state when an episode starts.

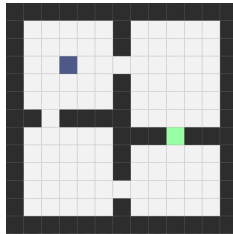


Figure 1: The Four Room environment. The black cells indicate the wall. The white cells indicate the area the agents can be in. The blue cell indicates the current position of the agent. The green cell indicates the goal.

3 Related Work

In this section, some option discovery algorithms will be summarized.

Option-Critic

Option-Critic [8] is an RL algorithm inspired by Action-Critic [5], where options are trained to maximize expected return, while an option-value function Q_Ω is trained to evaluate the decision the options.

Deliberation Cost

Since optimal policy can be achieved even without using options, if options are trained only to maximize expected return, they may degenerate and either terminate every steps or never terminate. Deliberation Cost [9] is a way to encourage longer option duration by punishing option switching.

Interest Option-Critic

The original Option-Critic assumes that options can be initiated everywhere, Interest Option-Critic [2] tries to remove this assumption by introducing interest functions $I : \mathcal{S} \times \Omega \rightarrow [0, 1]$ as a replacement for the initiation set. Experimental result shows that options learned by Interest Option-Critic is localized.

Termination-Critic

Termination-Critic [10] changes the objective of the termination function β from maximizing the expected return to minimize the entropy of the termination state. Since entropy can be interpreted as the information gain, this means minimizing the information gained from knowing the termination state, or in other words, making the termination state more predictable.

Attention Option-Critic

Attention Option-Critic [3] implements attention mechanism into Option-Critic. Different options are trained to attend to different features of the state. The attention units were trained to not only maximize the expected return, but also other things like maximizing difference between attention of different options.

4 Exploration

An analysis on localization will be conducted in this section.

What is Localization?

Localization is about options each responsible for a sub-task, or another way of looking at it is options each representing a skill. However, defining and measuring localization quantitatively is hard, which is why most work evaluate these option discovery algorithms qualitatively, by observing the agent acting for an episode in the environment. For example in the Four Rooms environment, only using one option in each room is considered as localization.

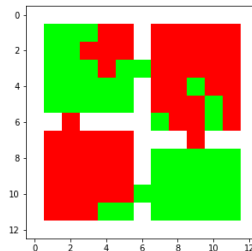


Figure 2: Red and green each represents an option, the options are pretty localized here.

Why Localization?

To understand why we want localization, first we need to answer a fundamental question: Why do we even use options in the first place? Is it to maximize expected return? However, optimal policy can

be achieved using only primitive actions. If options cannot give us a higher return, why do we even need options? Some researchers suggest that options should speed up planning [9] [10] and also options should be transferable [2] [3].

If this is what a good option should be like, then a set of localized options would be beneficial. Localized options are easy to interpret, which made it easily reusable when transferred to a different environment. Also, easy-to-interpret options can speed up planning because each options have its clear purpose and usage.

How Localization is achieved?

Now I will analyze how some of the previous work achieve localization of options.

Attention Option-Critic

In Attention Option-Critic [3], each options are trained to attend to different features of the state. My hypothesis is that the attention mechanism can act as a constraint on what kind of policy each option can have. Each features of the state represents a piece of information about the state. When performing a sub-task, not all the features are necessary. Each sub-task requires different subset of features. Since the attention mechanism limits the subset of features given to an option, the option cannot learn sub-task that requires features outside of the subset of features it was given, or else the option will perform poorly. In the algorithm, each option is trained to have diverse attention, which force each option to learn to complete a different sub-task.

For example, there is 3 options and an RGB 2D image is the features of the state. Suppose the 3 options each attend to one of the RGB channels, and one of the sub-tasks is checking if there is a purple circle on the image. In this case, only the option with attention on the green channel can complete this sub-task.

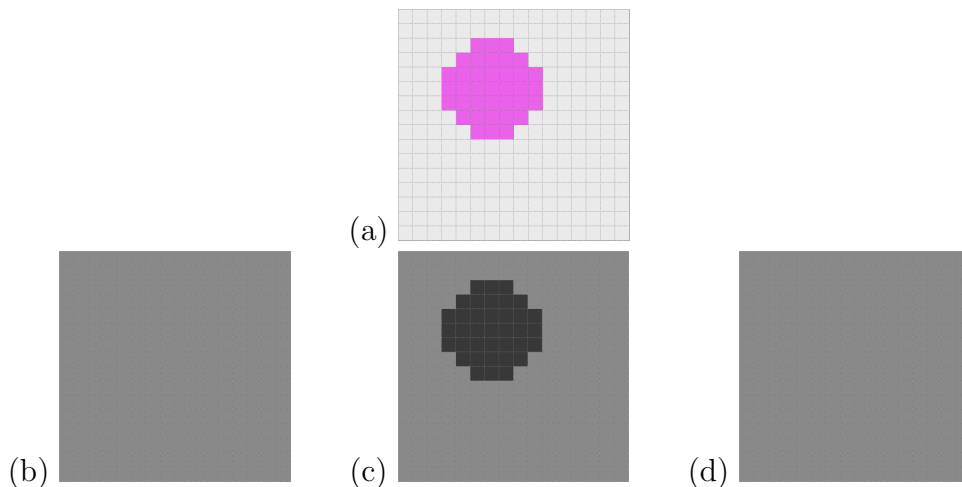


Figure 3: (a) is the RGB input, (b), (c) and (d) are the red, green and blue channel respectively. The circle can only be seen in the green channel.

Deliberation Cost

In Deliberation Cost [9], options are encouraged to be more temporally extended. Since options that terminates every step must be non-localized, Deliberation Cost can increase the chance of achieving localization.

Termination-Critic

In Termination-Critic [10], option termination states' entropy is being minimized, and experimental

results show that option trained by this usually choose to terminate in bottleneck states (frequently visited states). My hypothesis is that bottleneck states are usually the start or end of a sub-task, having the option terminate at these states essentially chains termination with initiation.

I will illustrate this with a simple example: In the Four Rooms environment, assume that the sub-task is walking from one doorway to another. The two doorway are bottleneck states because the agent must go through them. Since the agent can take on many paths, all the other states are not bottleneck states. When the agent get to the next doorway, another option can be immediately initiated.

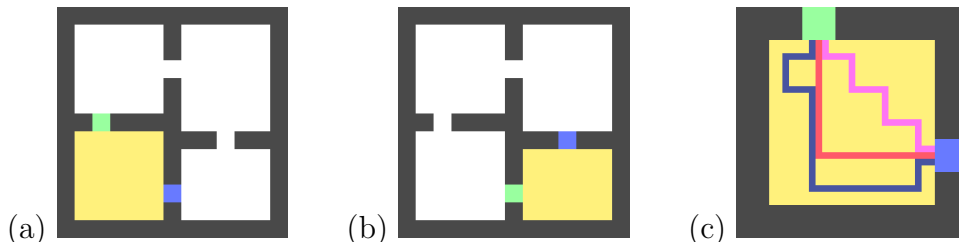


Figure 4: (a) and (b) are two options in the theoretic example. Green is the start of the option, blue is the end of the option, yellow is the intermediate states the option may encounter. (c) is a zoomed-in version of the bottom left room. Dark blue, red and pink lines are some of the paths the option can take.

Interest Option-Critic

In Interest Option-Critic. My hypothesis is that the interest function made the policy-over-options bias to choosing one of the options in a state, and since the paper uses a neural network as the interest function, the policy-over-options will also bias to choosing that option in a neighboring states.

5 The Localization Framework

Naturally, the next question that will be asked is: What do all of these algorithms have in common? Now I will propose a framework that can act as an abstraction for all of these algorithms.

The Localization Framework

1. Grouping

Group states into meaningful sub-tasks based on a certain criterion

2. Assignment

Assign the sub-tasks to different options

3. Optimization

Train options to perform well in the sub-task it is given and also improve the initial grouping of sub-tasks

4. Selection

Form the policy-over-options to select option in different state

This framework is inspired by Adaboost [11], which is an Ensemble Learning algorithm from Supervised Learning. There are a lot of similarities between Ensemble Learning and Option Learning, this has already been pointed out in previous work [12], the individual weak classifiers can be thought of as options.

Each weak classifier is responsible for classifying a small subset of the training data, just like how each option is responsible for a sub-task. In Adaboost, a bunch of weak classifiers are trained sequentially, each of them focuses on training data that is classified poorly by the previous weak classifiers.

Since this training process involves dividing training example into groups, then assign it to different weak classifiers, it inspires the Grouping and Assignment steps in the Localization Framework. Also, the

Optimization step in the framework is reminiscent of the weak classifiers learning to classify the training examples. After a lot of weak classifiers are trained, Adaboost combines them together to form a boosted classifier. The boosted classifier is the weighted sum of all the weak classifiers, the weighting is somewhat like a selection process, so it inspires the Selection step in the framework.

	Grouping	Assignment	Optimization	Selection
Attention Option-Critic	Group states that need the same sets of features	The algorithm assign an attention mechanism to each option	Options and attention mechanisms maximize return	Choose the option with maximum expected return
Termination-Critic	Group states between two bottleneck states	Each option terminates in a bottleneck state	Internal policy maximize return, termination function minimize entropy	Choose the option with maximum expected return
Interest Option-Critic	Group states that are close together	The algorithm assign an interest function to each option	Options and interest functions maximize return	Choose the option with high expected return and interest

Table 1: Previously mentioned algorithms can be fitted into the Localization Framework

This framework is the abstraction of algorithms that produce localized options, so it is very useful in deriving a new algorithm in the next section.

6 Attention-Over-Actions Option-Critic

Now that there is a framework, I can just follow the framework and derive a new algorithm. The following algorithm will be called Attention-Over-Actions Option-Critic because it perform abstraction on the action space, this algorithm is largely inspired by Attention Option-Critic.

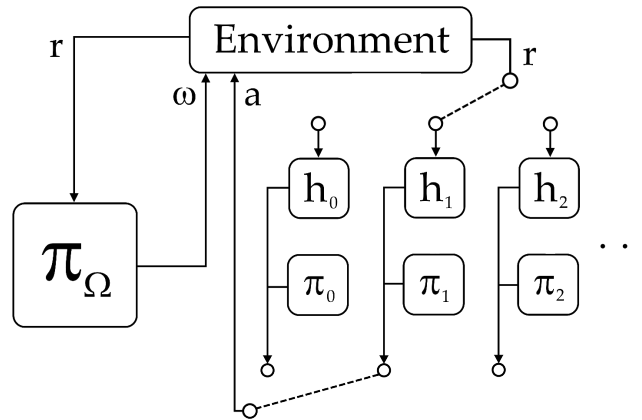


Figure 5: Visualization of the interaction between the environment and the Attention-Over-Actions Option-Critic algorithm.

1. Grouping

In Attention Option-Critic, the Grouping step divide sub-tasks based on features needed, this works because each sub-task requires a different subset of features. The following algorithm borrow the idea of attention over state features and apply it to actions, so each option will attend to a different sets of actions.

The intuition for this is that each sub-task will not need to use all the actions, for example, you would not consider performing a kicking action when you are driving a car. This is essentially performing actions abstraction, since the option can only attend to a subset of the actions.

2. Assignment

The following algorithm assigns the sub-task in a similar way as Attention Option-Critic, an attention mechanism $h_{\omega,\phi} : \mathcal{A} \rightarrow [0, 1]$ parameterized by ϕ will be given to each option. But instead of masking the state features, it masks the probability for choosing each actions. So the final probability π_{h_ω} for option ω to choose action a will be:

$$\pi_{h_\omega}(a|s) = \frac{\pi_{\omega,\theta}(a|s)h_{\omega,\phi}(a)}{\sum_{a'} \pi_{\omega,\theta}(a'|s)h_{\omega,\phi}(a')}$$

where $\pi_{\omega,\theta}$ is the internal policy of option ω and is parameterized by θ .

3. Optimization

The optimization of the attention mechanism and that of the option will be described separately.

Option Optimization

Let's first consider the optimization of the option. I will train the internal policy and termination function just like in Option-Critic. They will maximize the expected return by using gradient ascent:

$$\theta \leftarrow \theta + \alpha_\theta \nabla_\theta Q_\Omega(s, \omega)$$

$$\nu \leftarrow \nu + \alpha_\nu \nabla_\nu U_\Omega(s', \omega)$$

where θ and ν are the parameters for the internal policy and termination function respectively, $Q_\Omega(s, \omega)$ and $U_\Omega(s', \omega)$ are the expected return for choosing the option ω in state s and entering the state s' when using option ω respectively.

I can directly reuse the result from the Option-Critic paper:

$$\nabla_\theta Q_\Omega(s, \omega) = E[\nabla_\theta \log \pi_{\omega,\theta}(a|s) Q_U(s, \omega, a)]$$

$$\nabla_\nu U_\Omega(s', \omega) = E[\nabla_\nu \beta_{\omega,\nu}(s') A(s', \omega)]$$

where $Q_U(s, \omega, a)$ is the expected return for choosing action a when using option ω in state s .

The reason why I ignored the attention mechanism here is that the internal policy should not know the about the attention mechanism, or else it may want to revert the effect of the attention mechanism, i.e. Assigning a high weight to an action with low attention.

Attention Mechanism Optimization

Now let's consider the optimization of $h_{\omega,\phi}$. I want the grouping to be different for all options because or else all options will just aim for the sub-task with highest return. I also want the options to focus on as little actions as possible while still having acceptable performance. Essentially what I need is the algorithm to consider the trade-off between these objectives and achieve a balance between them. A nice way to do this is to add all of these objective up and then perform gradient ascent on the sum:

$$\phi \leftarrow \phi + \alpha_\phi \nabla_\phi \sum_o (w_o O_o)$$

where o is the index of an objective, w_o is the weight of the objective, O_o is the objective function. This method has been used for Attention Option-Critic too. Now I will list out the objectives that I want the option to consider:

1. Perform well
2. Different from other options
3. The components of the attention mechanism is close to 0 or 1
4. Focus on small set of actions

For the first objective, I can just use $Q_\Omega(s, \omega)$ like in Attention Critic.

$$\max_h O_1 = \max_h Q_\Omega(s, \omega)$$

For the second objective, I will minimize cosine similarity just like in Attention Critic.

$$\min_h O_2 = \min_h \sum_{h' \neq h} C(h, h') = \min_h \sum_{h' \neq h} \frac{\langle h', h \rangle}{\|h'\| \times \|h\|}$$

For the third objective, I will minimize entropy in the attention mechanism. Entropy measures the uncertainty in a probability distribution, so h should be normalized first.

$$\max_h O_3 = \max_h H\left(\frac{h}{\|h\|}\right) = \max_h \left\langle \frac{h}{\|h\|}, \log \frac{h}{\|h\|} \right\rangle$$

For the forth objective, I will minimize the length of the attention mechanism, which discourage focusing on too many actions.

$$\min_h O_4 = \min_h \|h\|$$

4. Selection

Any policy-over-options that favor higher Q-value options will work in this case, because the option will need the right set of actions in order to perform well, or else it will fail horribly. So the Q-value already encoded which option has the right set of actions. This means that policies like ϵ -greedy should work for this algorithm.

Algorithm 1 Pseudocode for Attention-Over-Actions Option-Critic (AOAOC)

```

 $s \leftarrow s_0$ 
Choose  $\omega$  according to the policy-over-options  $\pi_\Omega(s)$ 
repeat
    Choose  $a$  according to  $\pi_{h_\omega}(a|s)$ 
    Take action  $a$  in  $s$ , observe  $s', r$ 

    1. Options evaluation:
     $\delta \leftarrow r - Q_U(s, \omega, a)$ 
    if  $s'$  is non-terminal then
         $\delta \leftarrow \delta + \gamma(1 - \beta_{\omega, \nu}(s'))Q_\Omega(s', \omega) + \gamma\beta_{\omega, \nu}(s')\max_{\omega'} Q_\Omega(s', \omega')$ 
    end if
     $Q_U(s, \omega, a) \leftarrow Q_U(s, \omega, a) + \alpha\delta$ 

    1. Options improvement:
     $\theta \leftarrow \theta + \alpha_\theta \nabla_\theta \log \pi_{\omega, \theta}(a|s)Q_U(s, \omega, a)$ 
     $\nu \leftarrow \nu + \alpha_\nu \nabla_\nu \beta_{\omega, \nu}(s')(Q_\Omega(s', \omega) - V_\Omega(s'))$ 
     $\phi \leftarrow \phi + \alpha_\phi \nabla_\phi \sum_o (w_o O_o)$ 
    if  $\beta_{\omega, \nu}$  terminates in  $s'$  then
        choose new  $\omega$  according to the policy-over-options  $\pi_\Omega(s)$ 
    end if
     $s \leftarrow s'$ 
until  $s'$  is terminal

```

7 Experiment

In this section the algorithm AOAOC will be tested in the Four Rooms environment.

Initial Setup of the Experiment

Q_Ω and Q_U each is represented as a tensor. I do not need to use a neural network here since the state space is discrete in Four Rooms.

π_Ω is an ϵ -greedy policy for the Q_Ω , i.e. the policy-over-options has a probability of $1 - \epsilon + \frac{\epsilon}{n_\Omega}$ to select the best option, where n_Ω is the number of options.

ν is a tensor, and $\beta_{\omega,\nu}$ is parameterized by ν . More precisely, $\beta_{\omega,\nu} = \sigma(\nu)$, where σ is the sigmoid function $\frac{1}{1+e^{-x}}$. This can ensure that $\beta_{\omega,\nu}$ is a probability.

θ is also a tensor, and $\pi_{\omega,\theta}$ is parameterized by θ . More precisely, $\pi_{\omega,\theta}$ is the softmax over the components of different actions in θ . This can ensure that $\pi_{\omega,\theta}$ is a probability distribution over actions.

ϕ is also a tensor, and $h_{\omega,\phi}$ is parameterized by ϕ . Just like as $\beta_{\omega,\nu}$, $h_{\omega,\phi} = \sigma(\phi)$. This can ensure that $h_{\omega,\phi}$ is between 0 and 1.

Also, I would like the agent to have a larger incentive to go to the goal in a shorter duration, so I added a punishment of -2 for each step taken by the agent.

Problematic Attention Mechanism

When running the experiment, two problem is encountered:

1. The attention mechanism h_ω often becomes all one.
2. The attention objectives conflict with one another.

Problem 1 leads to degenerate attention mechanism because the option is paying the same amount of attention to all actions. Suppose each of the component in the attention mechanism is a constant k , i.e. $h_\omega = [k, k, k, \dots]$. The final policy will be:

$$\pi_{h_\omega}(a|s) = \frac{k \times \pi_\omega(a|s)}{\sum_{a'} k \times \pi_\omega(a'|s)} = \frac{\pi_\omega(a|s)}{\sum_{a'} \pi_\omega(a'|s)}$$

which is the same as not using attention at all.

This problem is probably caused by the objective 1, which is to maximizing $Q_\Omega(s, \omega)$. Having a uniform attention always yield a higher return because there will not be a constraint to which action it can use.

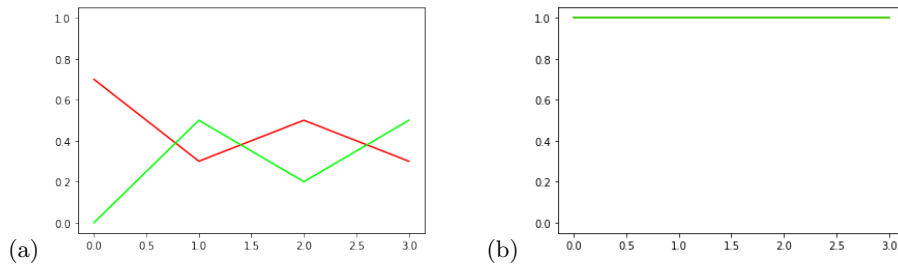


Figure 6: x-axis is the actions. y-axis is the attention on that action. Red and green lines each represents an attention mechanism. (a) is what the attention mechanism should look like. (b) is when it becomes all ones.

Problem 2 is follows directly from problem 1. Objective 1 will try to increase all the components to 1. Objective 3 will try to make only one of the components to 1, while the others all 0. Objective 4 will try to decrease all the components to 0.

My original intention is that the algorithm will automatically balance each objective and find the optimal attention mechanism. However, it turns out that the attention mechanism almost always result in 1 of 3 situations:

1. All components are 1

2. All components are 0
3. Only 1 of the component is one, while the others are all zero

For most of the time, 1 of the objectives completely dominates and results in 1 of the 3 situations above.

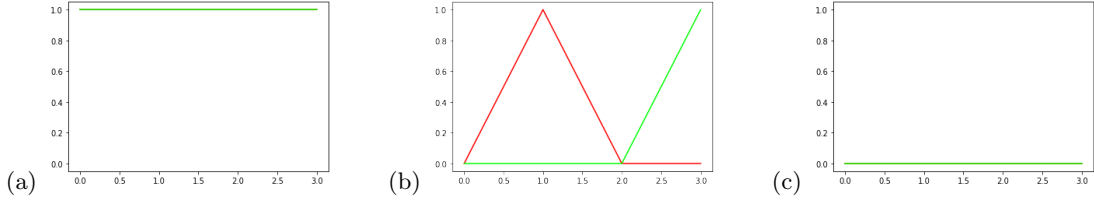


Figure 7: (a), (b) and (c) is the perfect end result for objective 1,3 and 4 respectively. All of which are degenerated attention mechanism.

8 Conclusion

9 Future Direction

This algorithm currently only focuses on discrete action space, a possible future direction might be extending this to continuous space. For vector actions, it might be good to attend to only some of the components of the vector, while ignoring other components by either randomly selecting values for them or keeping them the same as in the last state of the previous option. For scalar action, one direction is to use a one dimensional gaussian distribution as the attention mechanism, and multiply it with the original action distribution. However, some sort of trick may need to be deployed to speed up the process of normalization.

Bibliography

- [1] R. S. Sutton, D. Precup, and S. Singh, “Between mdps and semi-mdps: A framework for temporal abstraction in reinforcement learning,” *Artificial Intelligence*, vol. 112, no. 1, pp. 181 – 211, 1999.
- [2] K. Kheterpal, M. Klissarov, M. Chevalier-Boisvert, P.-L. Bacon, and D. Precup, “Options of interest: Temporal abstraction with interest functions,” 2020.
- [3] R. Chunduru and D. Precup, “Attention option-critic,” 2020.
- [4] R. Sutton and A. Barto, *Reinforcement Learning: An Introduction*. Adaptive Computation and Machine Learning series, MIT Press, 2018.
- [5] V. Konda and J. Tsitsiklis, “Actor-critic algorithms,” 2000.
- [6] S. Bradtke and M. Duff, “Reinforcement learning methods for continuous-time markov decision problems,” 12 1994.
- [7] R. Sutton, D. Precup, and S. Singh, “Intra-option learning about temporally abstract actions,” pp. 556–564, 01 1998.
- [8] P.-L. Bacon, J. Harb, and D. Precup, “The option-critic architecture,” 2016.
- [9] J. Harb, P.-L. Bacon, M. Klissarov, and D. Precup, “When waiting is not an option : Learning options with a deliberation cost,” 2017.
- [10] A. Harutyunyan, W. Dabney, D. Borsa, N. Heess, R. Munos, and D. Precup, “The termination critic,” 2019.

- [11] Y. Freund and R. E. Schapire, “A decision-theoretic generalization of on-line learning and an application to boosting,” in *Computational Learning Theory* (P. Vitányi, ed.), (Berlin, Heidelberg), pp. 23–37, Springer Berlin Heidelberg, 1995.
- [12] S. Zhang, H. Chen, and H. Yao, “Ace: An actor ensemble algorithm for continuous control with tree search,” 2018.
- [13] K. Khetarpal, “ioc repository.” <https://github.com/kkhetarpal/ioc>, 2020.

Appendix

A Notations

Markov Decision Process

\mathcal{S} — Set of states
 \mathcal{A} — Set of actions
 r — Reward function or reward
 γ — Discount factor
 P — Transition model
 π — Policy
 s — State
 a — Action
 s_T — Terminal state

Option Framework

Ω — Set of options
 $\pi_{\omega, \theta}$ — Internal policy of option ω
 β_{ω} — Termination probability of option ω
 \mathcal{I}_{ω} — Initiation set of option ω
 π_{Ω} — Policy-over-options
 Q_U — Expected return for choosing an action
 Q_{Ω} — Expected return for choosing an option
 A_{Ω} — Expected advantage for choosing an option
 V_{Ω} — Expected return in a state
 U_{Ω} — Expected return for arriving in a state
 θ — Parameter for internal policy
 ν — Parameter for termination probability
 ω — option

Attention-Over-Actions Option-Critic

$h_{\omega, \phi}$ — Attention Mechanism
 $\pi_{h_{\omega}}$ — Final probability
 ϕ — Parameter for attention mechanism
 α — Learning rate for Q_U
 α_{θ} — Learning rate for internal policy
 α_{ν} — Learning rate for termination probability
 α_{ϕ} — Learning rate for attention mechanism
 O_o — Objective function
 w_o — Weight of objective function
 o — Objective
 δ — One step Q-value error

Operations

∇ — Gradient
 $E[\]$ — Expected value

\leftarrow — Assignment
 $<, >$ — Dot product
 $||\cdot||$ — Length of a vector

B Proof

B.1 Derivative of Objective 1

Let h_ω^a be the element in the attention vector h_ω that corresponds to the action a , i.e. $h_\omega = [h_\omega^{a_0}, h_\omega^{a_1}, \dots]$

$$\begin{aligned}
 \frac{\partial Q_\Omega(s, \omega)}{\partial h_\omega^a} &= \sum_{s, \omega} \mu_\Omega(s, \omega | s_0, \omega_0) \sum_a \frac{\partial \pi_{h_\omega}(a|s)}{\partial h_\omega^a} Q_U(s, \omega, a) \\
 &= E_{s, \omega, a \sim \pi_{h_\omega}} \left[\frac{\partial \log \pi_{h_\omega}(a|s)}{\partial h_\omega^a} Q_U(s, \omega, a) \right]
 \end{aligned}$$

The steps above follows directly from the Option-Critic appendix. The next step is to express $\frac{\partial \log \pi_{h_\omega}(a|s)}{\partial h_\omega^a}$ in a simpler form.

$$\begin{aligned}
 \frac{\partial \log \pi_{h_\omega}(a|s)}{\partial h_\omega^a} &= \frac{\partial}{\partial h_\omega^a} \left[\log \frac{\pi_\omega(a|s) h_\omega^a}{\sum_{a'} \pi_\omega(a'|s) h_\omega^{a'}} \right] \\
 &= \frac{\partial}{\partial h_\omega^a} [\log \pi_\omega(a|s) + \log h_\omega^a - \log \sum_{a'} \pi_\omega(a'|s) h_\omega^{a'}] \\
 &= \frac{\partial}{\partial h_\omega^a} [\log \pi_\omega(a|s)] + \frac{\partial}{\partial h_\omega^a} [\log h_\omega^a] - \frac{\partial}{\partial h_\omega^a} [\log \sum_{a'} \pi_\omega(a'|s) h_\omega^{a'}]
 \end{aligned}$$

$\pi_\omega(a|s)$ is independent of h_ω^a , hence $\frac{\partial}{\partial h_\omega^a} [\log \pi_\omega(a|s)] = 0$

$$\begin{aligned}
 &= \frac{\partial}{\partial h_\omega^a} [\log h_\omega^a] - \frac{\partial}{\partial h_\omega^a} [\log \sum_{a'} \pi_\omega(a'|s) h_\omega^{a'}] \\
 &= \frac{1}{h_\omega^a} - \frac{\pi_\omega(a|s)}{\sum_{a'} \pi_\omega(a'|s) h_\omega^{a'}} \\
 &= \frac{1}{h_\omega^a} - \frac{\pi_\omega(a|s)}{\sum_{a'} \pi_\omega(a'|s) h_\omega^{a'}} \times \frac{h_\omega^a}{h_\omega^a}
 \end{aligned}$$

By definition $\pi_{h_\omega}(a|s) = \frac{\pi_\omega(a|s) h_\omega^a}{\sum_{a'} \pi_\omega(a'|s) h_\omega^{a'}}$

$$\begin{aligned}
 &= \frac{1}{h_\omega^a} - \frac{\pi_{h_\omega}(a|s)}{h_\omega^a} \\
 &= \frac{1 - \pi_{h_\omega}(a|s)}{h_\omega^a}
 \end{aligned}$$

Substitute back to the original calculation,

$$\frac{\partial Q_\Omega(s, \omega)}{\partial h_\omega^a} = E_{s, \omega, a \sim \pi_{h_\omega}} \left[\frac{1 - \pi_{h_\omega}(a|s)}{h_\omega^a} Q_U(s, \omega, a) \right]$$

B.2 Derivative of Objective 2

Let h_ω^a be the element in the attention vector h_ω that corresponds to the action a , i.e. $h_\omega = [h_\omega^{a_0}, h_\omega^{a_1}, \dots]$. Since we would like to minimize the cosine similarity, the negative cosine similarity will be used instead.

$$\begin{aligned} \sum_{h_{\omega'} \neq h_\omega} \frac{\partial C(h_\omega, h_{\omega'})}{\partial h_\omega^a} &= \sum_{h_{\omega'} \neq h_\omega} \frac{\partial}{\partial h_\omega^a} \frac{\langle h_\omega, h_{\omega'} \rangle}{\|h_\omega\| \times \|h_{\omega'}\|} \\ &= \sum_{h_{\omega'} \neq h_\omega} \frac{\|h_\omega\| \times \|h_{\omega'}\| \frac{\partial}{\partial h_\omega^a} \langle h_\omega, h_{\omega'} \rangle - \langle h_\omega, h_{\omega'} \rangle \frac{\partial}{\partial h_\omega^a} (\|h_\omega\| \times \|h_{\omega'}\|)}{(\|h_\omega\| \times \|h_{\omega'}\|)^2} \end{aligned}$$

I will calculate each of the derivative separately:

$$\begin{aligned} \frac{\partial}{\partial h_\omega^a} \langle h_\omega, h_{\omega'} \rangle &= \frac{\partial}{\partial h_\omega^a} \sum_{a'} h_\omega^{a'} h_{\omega'}^{a'} = h_{\omega'}^a \\ \frac{\partial}{\partial h_\omega^a} \|h_\omega\| \times \|h_{\omega'}\| &= \|h_{\omega'}\| \frac{\partial}{\partial h_\omega^a} \sqrt{\sum_{a'} (h_\omega^{a'})^2} = \|h_{\omega'}\| \frac{h_\omega^a}{\sqrt{\sum_{a'} (h_\omega^{a'})^2}} \end{aligned}$$

Substitute back to the original calculation,

$$\begin{aligned} \sum_{h_{\omega'} \neq h_\omega} \frac{\partial}{\partial h_\omega^a} \frac{\langle h_\omega, h_{\omega'} \rangle}{\|h_\omega\| \times \|h_{\omega'}\|} &= \sum_{h_{\omega'} \neq h_\omega} \frac{\|h_\omega\| \times \|h_{\omega'}\| h_{\omega'}^a - \langle h_\omega, h_{\omega'} \rangle \|h_{\omega'}\| \frac{h_\omega^a}{\sqrt{\sum_{a'} (h_\omega^{a'})^2}}}{(\|h_\omega\| \times \|h_{\omega'}\|)^2} \\ &= \sum_{h_{\omega'} \neq h_\omega} \frac{h_{\omega'}^a}{\|h_\omega\| \times \|h_{\omega'}\|} - \frac{\langle h_\omega, h_{\omega'} \rangle h_\omega^a}{(\|h_\omega\| \times \|h_{\omega'}\|) \times \|h_\omega\|^2} \end{aligned}$$

B.3 Derivative of Objective 3

Let h_ω^a be the element in the attention vector h_ω that corresponds to the action a , i.e. $h_\omega = [h_\omega^{a_0}, h_\omega^{a_1}, \dots]$. Since entropy usually applies to probabilities, I will normalize the attention unit into $\bar{h}_\omega = \frac{h_\omega}{\sum_{a'} h_\omega^{a'}}$

$$\begin{aligned} \frac{\partial H(\bar{h}_\omega)}{\partial h_\omega^a} &= \frac{\partial \sum_{a'} \bar{h}_\omega^{a'} \log \bar{h}_\omega^{a'}}{\partial h_\omega^a} \\ &= \frac{\partial \bar{h}_\omega^a \log \bar{h}_\omega^a}{\partial h_\omega^a} + \sum_{\bar{h}_\omega^{a'} \neq \bar{h}_\omega^a} \frac{\partial \bar{h}_\omega^{a'} \log \bar{h}_\omega^{a'}}{\partial h_\omega^a} \\ &= \frac{\partial \bar{h}_\omega^a \log \bar{h}_\omega^a}{\partial \bar{h}_\omega^a} \times \frac{\partial \bar{h}_\omega^a}{\partial h_\omega^a} + \sum_{\bar{h}_\omega^{a'} \neq \bar{h}_\omega^a} \frac{\partial \bar{h}_\omega^{a'} \log \bar{h}_\omega^{a'}}{\partial \bar{h}_\omega^{a'}} \times \frac{\partial \bar{h}_\omega^{a'}}{\partial h_\omega^a} \end{aligned}$$

I will calculate each of the derivative separately:

$$\begin{aligned} \frac{\partial \bar{h}_\omega^a \log \bar{h}_\omega^a}{\partial \bar{h}_\omega^a} &= \log \bar{h}_\omega^a + \bar{h}_\omega^a \frac{\partial \log \bar{h}_\omega^a}{\partial \bar{h}_\omega^a} = \log \bar{h}_\omega^a + 1 \\ \frac{\partial \bar{h}_\omega^a}{\partial h_\omega^a} &= \frac{\partial}{\partial h_\omega^a} \frac{h_\omega^a}{\sum_b h_\omega^b} = \frac{\sum_b h_\omega^b - h_\omega^a}{(\sum_b h_\omega^b)^2} \\ \frac{\partial \bar{h}_\omega^{a'}}{\partial h_\omega^a} &= \frac{\partial}{\partial h_\omega^a} \frac{h_\omega^{a'}}{\sum_b h_\omega^b} = \frac{-h_\omega^{a'}}{(\sum_b h_\omega^b)^2} \end{aligned}$$

Substitute back to the original calculation,

$$\begin{aligned}\frac{\partial \bar{h}_\omega^a \log \bar{h}_\omega^a}{\partial \bar{h}_\omega^a} &= (\log \bar{h}_\omega^a + 1) \frac{\sum_b h_\omega^b - h_\omega^a}{(\sum_b h_\omega^b)^2} + \sum_{h_\omega^{a'} \neq h_\omega^a} (\log \bar{h}_\omega^{a'} + 1) \frac{-h_\omega^{a'}}{(\sum_b h_\omega^b)^2} \\ &= \frac{(\log \bar{h}_\omega^a + 1)}{\sum_b h_\omega^b} + \sum_{h_\omega^{a'}} (\log \bar{h}_\omega^{a'} + 1) \frac{-h_\omega^{a'}}{(\sum_b h_\omega^b)^2}\end{aligned}$$

B.4 Derivative of Objective 4

Let h_ω^a be the element in the attention vector h_ω that corresponds to the action a , i.e. $h_\omega = [h_\omega^{a_0}, h_\omega^{a_1}, \dots]$. Since we would like to minimize the length, the negative length will be used instead.

$$\begin{aligned}\frac{\partial ||h_\omega||}{\partial h_\omega^a} &= \frac{\partial}{\partial h_\omega^a} \sqrt{\sum_{a'} (h_\omega^{a'})^2} \\ &= \frac{h_\omega^a}{\sqrt{\sum_{a'} (h_\omega^{a'})^2}} = \frac{h_\omega^a}{||h_\omega||}\end{aligned}$$

C Experimental Details

D Code

The code below is based on codes in the ioc repository.[13]

fourrooms.py

```

1 import numpy as np
2 from random import uniform
3
4 class Fourrooms:
5     def __init__(self, initstate_seed, punishEachStep, deterministic, modified, easier):
6         self.punishEachStep = punishEachStep
7         self.deterministic = deterministic
8         self.modified = modified
9         if easier:
10             self.layout = """\
11 wwwwww
12 w      w
13 w      w
14 w      w
15 w      w
16 w      w
17 w wwwww
18 w      wwwww
19 w      w
20 w      w
21 w      w
22 w      w
23 wwwwww
24 """
25         else:
26             self.layout = """\
27 wwwwww
28 w      w
29 w      w
30 w      w
31 w      w

```

```

32 W      W      W
33 WW  WWW      W
34 W      WWW  WWW
35 W      W      W
36 W      W      W
37 W      W      W
38 W      W      W
39 WWWWWWWWWWWW
40 ""
41
42
43     self.occupancy = np.array([list(map(lambda c: 1 if c=='w' else 0, line)) for line in
44                               self.layout.splitlines()])
45
46     self.action_space = 4
47
48     self.observation_space = int(np.sum(self.occupancy == 0))
49
50     # 0 - Up
51     # 1 - Down
52     # 2 - Left
53     # 3 - Right
54
55     self.directions = [np.array((-1,0)), np.array((1,0)), np.array((0,-1)), np.array
56                       ((0,1))]
57
58     self.rng = np.random.RandomState(1234)
59
60     self.initstate_seed = initstate_seed
61     self.rng_init_state = np.random.RandomState(self.initstate_seed)
62
63     self.tostate = {}
64
65     self.occ_dict = dict(zip(range(self.observation_space),
66                             np.argwhere(self.occupancy.flatten() == 0).squeeze()))
67
68     statenum = 0
69     for i in range(13):
70         for j in range(13):
71             if self.occupancy[i, j] == 0:
72                 self.tostate[(i, j)] = statenum
73                 statenum += 1
74
75     self.tocell = {v:k for k,v in self.tostate.items()}
76
77     self.goal = 62
78     self.init_states = list(range(self.observation_space))
79     self.init_states.remove(self.goal)
80
81     def empty_around(self, cell):
82         avail = []
83         for action in range(self.action_space):
84             nextcell = tuple(cell + np.multiply(self.directions[action], self.inQuad24(self.
85             currentcell)))
86             if not self.occupancy[nextcell]:
87                 avail.append(nextcell)
88         return avail
89
90     def reset(self, test=None):
91         if test:
92             state = test
93         else:
94             state = self.rng_init_state.choice(self.init_states)
95         self.currentcell = self.tocell[state]
96         return state
97
98     def step(self, action):

```

```

97     reward = -2 * int(self.punishEachStep)
98     if self.rng.uniform() < 1/3 and not(self.deterministic):
99         empty_cells = self.empty_around(self.currentcell)
100         nextcell = empty_cells[self.rng.randint(len(empty_cells))]
101     else:
102         nextcell = tuple(self.currentcell + np.multiply(self.directions[action], self.
103             inQuad24(self.currentcell)))
104
105     if not self.occupancy[nextcell]:
106         self.currentcell = nextcell
107
108     state = self.tostate[self.currentcell]
109
110     if state == self.goal:
111         reward = 50
112
113     done = state == self.goal
114     return state, reward, float(done), None
115
116 def inQuad24(self, cell):
117     if not(self.modified):
118         return np.array([1,1])
119     if cell[1] > 6:
120         if cell[0] < 7:
121             return np.array([1, -1])
122     else:
123         if cell[0] > 6:
124             return np.array([1, -1])
125     return np.array([1, 1])

```

aoaoc_tabular.py

```

1 import numpy as np
2 from fourrooms import Fourrooms
3 from scipy.special import logsumexp, expit, softmax
4 '''
5 =====CLASS MAP=====
6 Option
7     - FinalPolicy pi_h
8     - Internal Policy (SoftmaxPolicy) pi_omega
9     - Attention Unit (LearnableAttention/PredefinedAttention) h_omega
10         - Value Objective (ValueObj) o1
11         - Cosine Similarity Objective (CoSimObj) o2
12         - Entropy Objective (EntropyObj) o3
13         - Length Objective (LengthObj) o4
14     - Termination Function (SigmoidTermination) beta_omega
15     - Q_omega (Q_U)
16 Policy Over Options (POO) pi_Omega
17     - Policy (EgreedyPolicy)
18     - Q_Omega (Q_U)
19 '''
20 #=====Option=====
21 class Option:
22     def __init__(self, rng, nfeatures, nactions, args, policy_over_options, index):
23         self.weights = np.zeros((nfeatures, nactions))
24         self.policy = FinalPolicy(rng, nfeatures, nactions, args, self.weights, index)
25         self.termination = SigmoidTermination(rng, nfeatures, args)
26         self.Qval = Q_U(nfeatures, nactions, args, self.weights, policy_over_options)
27
28     def sample(self, phi):
29         return self.policy.sample(phi)
30
31     def terminate(self, phi, value=False):
32         if value:
33             return self.termination.pmf(phi)

```



```

34         else:
35             return self.termination.sample(phi)
36
37     def _Q_update(self, traject, reward, done, termination):
38         self.Qval.update(traject, reward, done, termination)
39
40     def _H_update(self, traject):
41         qVal = self.Qval.value(traject[0][0], traject[2])
42         self.policy.H_update(traject, qVal)
43
44     def _B_update(self, phi, option, advantage):
45         self.termination.update(phi, option, advantage)
46
47     def _P_update(self, traject, baseline):
48         self.policy.P_update(traject, baseline)
49
50     def update(self, traject, reward, done, phi, option, termination, baseline, advantage):
51         self._Q_update(traject, reward, done, termination)
52         self._H_update(traject)
53         self._P_update(traject, baseline)
54         self._B_update(phi, option, advantage)
55
56
57     #=====Final Policy=====
58     class FinalPolicy:
59         def __init__(self, rng, nfeatures, nactions, args, qWeight, index):
60             self.rng = rng
61             self.nactions = nactions
62             self.internalPI = SoftmaxPolicy(rng, nfeatures, nactions, args, qWeight)
63             if (args.h_learn):
64                 self.attention = LearnableAttention(nactions, args, index)
65             else:
66                 self.attention = PredefinedAttention(args, index)
67
68         def pmf(self, phi):
69             pi = self.internalPI.pmf(phi)
70             h = self.attention.pmf()
71             normalizer = np.dot(pi, h)
72             return (pi*h)/normalizer
73
74         def sample(self, phi):
75             return int(self.rng.choice(self.nactions, p=self.pmf(phi)))
76
77         def H_update(self, traject, qVal):
78             self.attention.update(traject, self.pmf(traject[0][0]), qVal)
79
80         def P_update(self, traject, baseline):
81             self.internalPI.update(traject, baseline)
82
83
84     #=====Internal Policy=====
85     class SoftmaxPolicy:
86         def __init__(self, rng, nfeatures, nactions, args, qWeight):
87             self.rng = rng
88             self.nactions = nactions
89             self.temp = args.temp
90             self.weights = np.zeros((nfeatures, nactions))
91             self.qWeight = qWeight
92             self.lr = args.lr_intra
93
94         def _value(self, phi, action=None):
95             if action is None:
96                 return np.sum(self.weights[phi, :], axis=0)
97             return np.sum(self.weights[phi, action], axis=0)
98
99         def pmf(self, phi):
100             v = self._value(phi)/self.temp
101             return np.exp(v - logsumexp(v))

```

```

102
103     def sample(self, phi):
104         return int(self.rng.choice(self.nactions, p=self.pmf(phi)))
105
106     def update(self, traject, baseline):
107         actions_pmf = self.pmf(traject[0][0])
108         critic = self.qWeight[traject[0][0], traject[2]]
109         if baseline:
110             critic -= baseline
111         self.weights[traject[0][0], :] -= self.lr*critic*actions_pmf
112         self.weights[traject[0][0], traject[2]] += self.lr*critic
113
114
115 #=====Attention=====
116 class LearnableAttention():
117     def __init__(self, nactions, args, index):
118         self.weights = np.random.uniform(low=-1, high=1, size=(nactions,))
119         self.lr = args.lr_attend
120         self.o1 = ValueObj(args)
121         self.o2 = CoSimObj(args, index)
122         self.o3 = EntropyObj(args)
123         self.o4 = LengthObj(args)
124         CoSimObj.add2list(self)
125         self.normalize = args.normalize
126
127     def pmf(self):
128         if self.normalize:
129             return np.clip softmax(self.weights), 0.05, None)
130         return expit(self.weights)
131
132     def _grad(self):
133         attend = self.pmf()
134         return attend*(1. - attend)
135
136     def attention(self, a):
137         return self.pmf()[a]
138
139     def update(self, traject, finalPmf, qVal):
140         hPmf = self.pmf()
141         gradList = [self.o1.grad(traject[0][0], traject[2], hPmf, finalPmf, qVal), self.o2.
142                     grad(hPmf), self.o3.grad(hPmf), self.o4.grad(hPmf)]
143         self.weights += self.lr * np.sum(gradList, axis=0) * self._grad()
144         if self.normalize:
145             self.normalizing()
146
147     def normalizing(self):
148         self.weights -= np.mean(self.weights)
149
150 class PredefinedAttention():
151     def __init__(self, args, index):
152         if (index==0):
153             self.weights = np.array([1, 1, 1, 1])
154         if (index==1):
155             self.weights = np.array([1, 1, 1, 1])
156
157     def pmf(self):
158         return self.weights
159
160     def attention(self, a):
161         return self.pmf()[a]
162
163     def update(self, traject, finalPmf, qVal):
164         pass
165
166
167 #=====Objectives=====
168 class Objective:

```

```

169     def __init__(self, weight):
170         self.weight = weight
171
172     def grad(self):
173         return None
174
175     def loss(self):
176         return None
177
178
179 class ValueObj(Objective):
180     def __init__(self, args):
181         super().__init__(args.wo1)
182
183     def grad(self, phi, a, hPmf, finalPmf, qVal):
184         return self.weight * ((finalPmf + 1)/hPmf[a]) * qVal
185
186     def loss(self):
187         pass
188
189
190 class CoSimObj(Objective):
191     hList = []
192
193     def __init__(self, args, index):
194         super().__init__(args.wo2)
195         self.index = index
196
197     def grad(self, hPmf):
198         gradient = []
199         for i in range(len(hPmf)):
200             derivative = 0.
201             exclude = 0
202             for a in self.hList:
203                 if exclude == self.index:
204                     continue
205                 exclude += 1
206
207                 normalizer = np.linalg.norm(hPmf)*np.linalg.norm(a.pmf())
208                 term1 = a.pmf()[i]/normalizer
209                 term2 = hPmf[i]*np.dot(hPmf,a.pmf()) / (normalizer*np.power(np.linalg.norm(
210                     hPmf),2))
211                 derivative += -1*(term1 - term2)
212             gradient.append(derivative)
213         return self.weight * np.array(gradient)
214
215     def loss(self):
216         return np.sum([np.dot(hPmf,a.pmf())/(np.linalg.norm(hPmf)*np.linalg.norm(a.pmf()))
217             for a in self.hList])
218
219     @classmethod
220     def add2list(cls, attention):
221         cls.hList.append(attention)
222
223     @classmethod
224     def reset(cls):
225         cls.hList = []
226
227
228 class EntropyObj(Objective):
229     def __init__(self, args):
230         super().__init__(args.wo3)
231
232     def grad(self, hPmf):
233         gradient = []
234         normalizer = np.sum(hPmf)
235         normh = hPmf/normalizer
236         for i in range(len(hPmf)):

```

```

235         term1 = (1.+np.log(normh[i]))/normalizer
236         term2 = np.sum([(1.+np.log(normh[index]))*hPmf[index]/(normalizer**2) for index
237             in range(len(hPmf))])
238         gradient.append((term1-term2)*(self.loss(hPmf)-0.69))
239         return self.weight * np.array(gradient)
240
241     def loss(self, hPmf):
242         normalizer = np.sum(hPmf)
243         normh = hPmf/normalizer
244         return -1*np.sum(normh * np.log(normh))
245
246     class LengthObj(Objective):
247         def __init__(self, args):
248             super().__init__(args.wo4)
249
250         def grad(self, hPmf):
251             return -1 * hPmf / self.loss(hPmf)
252
253
254         def loss(self, hPmf):
255             return np.linalg.norm(hPmf)
256
257
258     =====Termination Function=====
259     class SigmoidTermination:
260         def __init__(self, rng, nfeatures, args):
261             self.rng = rng
262             self.weights = np.zeros((nfeatures,))
263             self.lr = args.lr_term
264             self.dc = args.dc
265
266         def pmf(self, phi):
267             return expit(np.sum(self.weights[phi]))
268
269         def sample(self, phi):
270             return int(self.rng.uniform() < self.pmf(phi))
271
272         def _grad(self, phi):
273             terminate = self.pmf(phi)
274             return terminate*(1. - terminate), phi
275
276         def update(self, phi, option, advantage):
277             magnitude, direction = self._grad(phi)
278             self.weights[direction] -= self.lr*magnitude*(advantage+self.dc)
279
280
281     =====Q-Value Individual Option=====
282     class Q_U:
283         def __init__(self, nfeatures, nactions, args, weights, policy_over_options):
284             self.weights = weights
285             self.lr = args.lr_criticA
286             self.discount = args.discount
287             self.policy_over_options = policy_over_options
288
289         def value(self, phi, action):
290             return np.sum(self.weights[phi, action], axis=0)
291
292         def update(self, traject, reward, done, termination):
293             update_target = reward
294             if not done:
295                 current_values = self.policy_over_options.value(traject[1][0])
296                 update_target += self.discount*((1. - termination)*current_values[traject[0][1]]
297                     + termination*np.max(current_values))
298
299             tderror = update_target - self.value(traject[0][0], traject[2])
300             self.weights[traject[0][0], traject[2]] += self.lr*tderror

```

```

301
302 #=====Policy Over Option=====
303 class P00:
304     def __init__(self, rng, nfeatures, args):
305         self.weights = np.zeros((nfeatures, args.noptions))
306         self.policy = EgreedyPolicy(rng, args, self.weights)
307         self.Q_Omega = Q_0(args, self.weights)
308
309     def update(self, traject, reward, done, termination):
310         self.Q_Omega.update(traject, reward, done, termination)
311
312     def sample(self, phi):
313         return self.policy.sample(phi)
314
315     def advantage(self, phi, option=None):
316         values = np.sum(self.weights[phi], axis=0)
317         advantages = values - np.max(values)
318         if option is None:
319             return advantages
320         return advantages[option]
321
322     def value(self, phi, option=None):
323         if option is None:
324             return np.sum(self.weights[phi, :], axis=0)
325         return np.sum(self.weights[phi, option], axis=0)
326
327 class EgreedyPolicy:
328     def __init__(self, rng, args, weights):
329         self.rng = rng
330         self.epsilon = args.epsilon
331         self.noptions = args.noptions
332         self.weights = weights
333
334     def _value(self, phi, action=None):
335         if action is None:
336             return np.sum(self.weights[phi, :], axis=0)
337         return np.sum(self.weights[phi, action], axis=0)
338
339     def sample(self, phi):
340         if self.rng.uniform() < self.epsilon:
341             return int(self.rng.randint(self.weights.shape[1]))
342         return int(np.argmax(self._value(phi)))
343
344
345 #=====Q-Value All Option=====
346 class Q_0:
347     def __init__(self, args, weights):
348         self.weights = weights
349         self.lr = args.lr_critic
350         self.discount = args.discount
351
352     def _value(self, phi, option=None):
353         if option is None:
354             return np.sum(self.weights[phi, :], axis=0)
355         return np.sum(self.weights[phi, option], axis=0)
356
357     def update(self, traject, reward, done, termination):
358         update_target = reward
359         if not done:
360             current_values = self._value(traject[1][0])
361             update_target += self.discount*((1. - termination)*current_values[traject[0][1]]
362                 + termination*np.max(current_values))
363
364         tderror = update_target - self._value(traject[0][0], traject[0][1])
365         self.weights[traject[0][0], traject[0][1]] += self.lr*tderror
366
367

```

```

368 #====Standard====
369 # Follow the code standard of the ioc repository
370 class Tabular:
371     def __init__(self, nstates):
372         self.nstates = nstates
373
374     def __call__(self, state):
375         return np.array([state,])
376
377     def __len__(self):
378         return self.nstates

```

visualize.py

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 from fourrooms import Fourrooms
4 from time import sleep
5
6 # 0 - Red
7 # 1 - Green
8 # 2 - Blue
9 # 3 - Black
10 class Visualization:
11     def __init__(self, fRoom, args, nactions, colorList
12                 =[[255,0,0],[0,255,0],[0,0,255],[0,0,0]]):
13         assert args.noptions <= len(colorList), "Length_of_color_list_must_match_number_of_options"
14         self.colorList = colorList
15         self.layout = fRoom.layout
16         self.occupancy = fRoom.occupancy
17         self.tostate = fRoom.tostate
18         self.tocell = fRoom.tocell
19         self.screen = np.array([list(map(lambda c: [0,0,0] if c=='w' else [255,255,255],
20                                         line)) for line in self.layout.splitlines()])
21         self.lastphi = None
22         self.noptions = args.noptions
23         self.nactions = nactions
24
25     def showMap(self, phi, option):
26         color = self.colorList[option]
27         self._draw(self.lastphi, [255,255,255])
28         self._draw(phi, color)
29         self.lastphi = phi
30         plt.figure(figsize=(5,5))
31         plt.subplot(111)
32         plt.imshow(self.screen, vmax=255, vmin=0)
33         plt.show()
34         sleep(0.05)
35
36     def showAttention(self, options):
37         x = np.array([i for i in range(self.nactions)])
38         plt.plot(x, np.array([int(i != 0) for i in range(self.nactions)]), color=[1,1,1])
39         for i in range(self.noptions):
40             plt.plot(x, options[i].policy.attention.pmf(), color=np.array(self.colorList[i])/255.)
41         plt.show()
42
43     def showPref(self, weight): # policy_over_options.weightsP or options[index].weightsP
44         for weight
45         pref = np.zeros((13,13,3), dtype="int")
46         for i in range(13):
47             for j in range(13):
48                 if self.occupancy[i,j] == 0:
49                     choice = np.argmax(weight[self.tostate[(i,j)],:])
50                     pref[i,j] = np.array(self.colorList[choice])

```

```

48         else:
49             pref[i,j] = np.array([255,255,255])
50     plt.figure(figsize=(5,5))
51     plt.subplot(111)
52     plt.imshow(pref, vmax=255, vmin=0)
53     plt.show()
54
55     def savePref(self, weight, algo=None, wo=None, dc=None, run=None):
56         pref = np.zeros((13,13,3), dtype="int")
57         for i in range(13):
58             for j in range(13):
59                 if self.occupancy[i,j] == 0:
60                     choice = np.argmax(weight[self.tostate[(i,j)],:])
61                     pref[i,j] = np.array(self.colorList[choice])
62                 else:
63                     pref[i,j] = np.array([255,255,255])
64     plt.figure(figsize=(5,5))
65     plt.subplot(111)
66     plt.imshow(pref, vmax=255, vmin=0)
67     if wo != None:
68         plt.savefig("../result/{0}/wo_{1}/dc_{2}/run_{3}.png".format(algo, str(wo), str(
69             dc), str(run)))
70     else:
71         plt.savefig("../result/{0}/dc_{1}/run_{2}.png".format(algo, str(dc), str(run)))
72
73     def resetMap(self, phi):
74         self.screen = np.array([list(map(lambda c: [0,0,0] if c=='w' else [255,255,255],
75             line)) for line in self.layout.splitlines()])
76         self.lastphi = phi
77         self._draw([62],[200,200,200])
78
79     def _draw(self, phi, rgb):
80         self.screen[self.tocell[phi[0]]] = np.array(rgb)

```