

Application of Modular Reinforcement Learning

Shun Zhang *

January 21, 2015

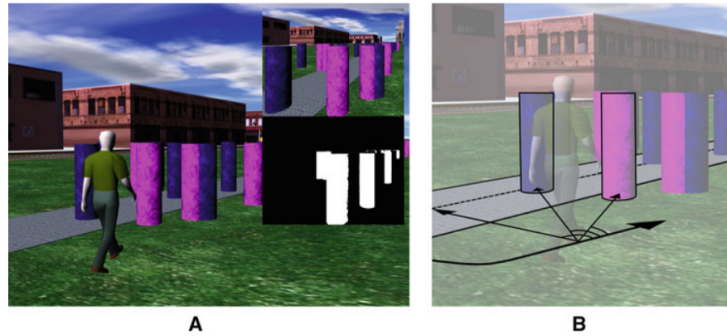


Figure 1:

Consider the task illustrated in Figure 1 A). The avatar is asked to do three sub-tasks simultaneously — 1) following the path, indicated by the gray line on the ground, 2) getting targets, the blue cylinders, and 3) avoiding obstacles, the pink cylinders [1].

From the reinforcement learning perspective, this task can be decomposed to be three sub-tasks as described above. In Figure 1 B), if the agent knows the distance and angle to an object, he is expected to know the optimal action to avoid or pursue it.

A human has not done this kind of task before can achieve a good performance easily, shown in the experiments described later. If a human knows the policies of the sub-tasks, or sub-MDPs, he can accomplish a complicated behavior by combining the sub-MDPs. That is,

$$Q(s, a) = \sum_i w_i Q_i(s, a)$$

where Q_i is the Q value of the i-th sub-MDP, w_i is the weight of the i-th sub-MDP. $w_i \geq 0, \sum_i w_i = 1$.

*This is part of the draft of my master thesis that I am working on in this semester (Spring 2015), only for the purpose of preview and evaluation of my thesis work. This is advised by Prof. Dana Ballard and Prof. Mary Hayhoe.

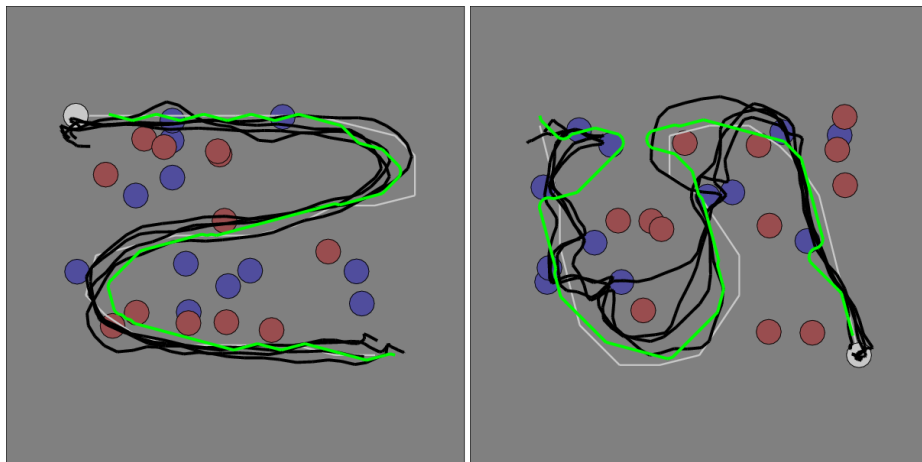


Figure 2:

Different weights can yield different performance. Let w_1, w_2, w_3 be weights for the task of target collection, obstacle avoidance, and path following, respectively. Let w be the vector of $(w_i)_1^n$. An agent with $w = (1, 0, 0)$ only collect targets, and one with $w = (0, 0.5, 0.5)$ may avoid the obstacles and follow the path.

From a different perspective, can we find a weight vector to best interpret human’s behavior? In Figure 2, there are two scenarios. Same as Figure 1, the red circles are obstacles. The blue circles are targets. The gray line is the path. The black lines are trajectories of human data. The left figure shows the case where human follows the path and ignores the targets and obstacles. The right figure shows the case that human does three sub-tasks simultaneously. The human data are collected by the Center for Perceptual Systems in University of Texas at Austin.

Now, we assume a learning agent that only knows the three sub-MDPs and the human data. It looks at the human behavior, and finds the weights that can interpret such behavior. Using such weights, the trajectories of our agents are drawn in the green lines. We can tell that in the left figure, the agent puts a large weight on path-following. In the right figure, it puts weights on all sub-MDPs.

References

- [1] Constantin A Rothkopf and Dana H Ballard. Modular inverse reinforcement learning for visuomotor behavior. *Biological cybernetics*, 107(4):477–490, 2013.