

UPF Report

Damien Allonsius & Anders Jonsson

30 April 2019

1 Introduction

2 Method

3 Initial results

4 Conclusion

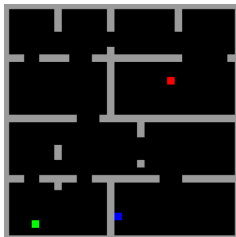
1 Introduction

2 Method

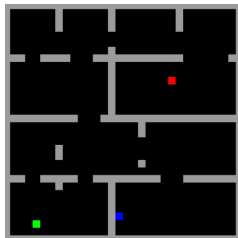
3 Initial results

4 Conclusion

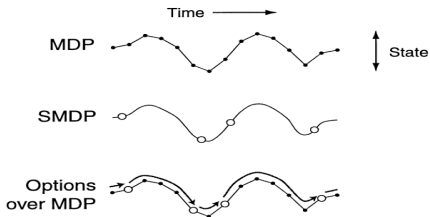
- An agent learns by experience to perform actions that achieve a certain task
- Already challenging for individual tasks
- Example : navigation task where the task is to pick up a key to open a door



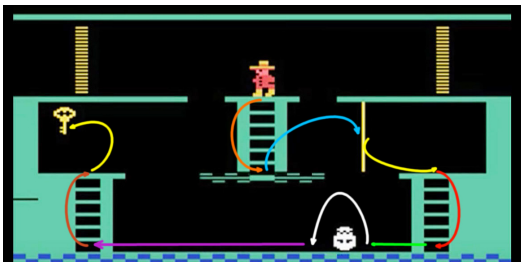
- No initial model of the environment
- Only available information is what agent observes right now
- Agent may face new tasks later
- The environment dynamics may change over time



- Decompose task into subtasks, each with its own policy
- Actions are performed at different time scales



Example : Montezuma's Revenge



- Exploration and planning at different levels of abstraction
- Capture subtasks that are relevant for a range of tasks
- Planning horizon is shorter for subtasks (easier to solve)
- Drawback : more learning problems !

1 Introduction

2 Method

3 Initial results

4 Conclusion

- Assume that the state space is partitioned into regions
(Formally, need a function $\Phi : S \rightarrow \mathbb{N}$ mapping states to partition IDs)
- Initially the agent is in a single region and has no options
- Two phases :
 - ① Exploration : find transitions among regions
 - ② Learning : introduce options that perform transitions between regions
- Abstract MDP : states are regions, actions are options

Example : Montezuma's Revenge

Partition : downsampled images of game screen
Option policies retain more detailed states

Normal view

Options' view

Agent's View

- At the agent level :
 - ① Exploration strategy with tree search.
 - ② When a reward is found : Q-Learning.

- At the option level, correct transitions can be learned with
 - ① Q-Learning
 - ② DQN (ongoing implementation).

- A reward (resp. penalty) is given when the option performs good (resp. bad) transitions decided by the agent.

- Can be combined with planning since agent maintains a transition model
- Flexible, can be combined with any strategy for identifying subtasks
- Partitions can be formed even when there is no clear task structure
- Transition options can potentially be reused (transfer learning)
- Options are only applicable in a single region (action elimination)

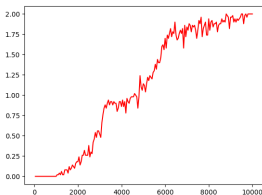
1 Introduction

2 Method

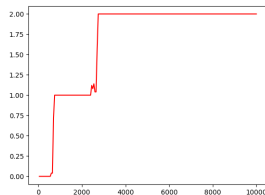
3 Initial results

4 Conclusion

Comparison between our strategy and Q-Learning. GRIDWORLD ENVIRONMENT



Q Learning



Hierarchy

MONTEZUMA'S REVENGE (ATARI)

We restart every time the agent picks the key.

With our strategy we get the key 19 times in 3000 iterations.

1 Introduction

2 Method

3 Initial results

4 Conclusion

Apply the hierarchical approach to the microgrid simulator :

- ① Downsample the state space by discretizing the state of charge of batteries (and the consumption and production of energy).
- ② Make options to go from a level of charge to another.

Let's talk about this together !