

**ECSE 526**  
**Assignment 3**  
Reinforcement Learning

Sean Stappas  
260639512

November 7<sup>th</sup>, 2017

## Introduction

### 1 Description of approach to generalization

Q-learning is the learning algorithm chosen for the agent. Many different distance metrics were used in the program, and close states are grouped based on this distance, as described by Mahadevan and Connell [4]. Here are three simple distance metrics, with the rationale for each choice of components.

**Manhattan** The Manhattan distance represents how close two states are on the Qbert pyramid, i.e., the number of blocks between them. This is a natural representation of distance in this environment, since Qbert can only change state by switching blocks and can only travel one block at a time (with the exception of using disks to jump back to the starting position).

**Hamming** Hamming distance represents the number of different state bits between two given states. The rationale behind this distance metric is that nearby states should have a similar bit pattern, assuming a sane state representation.

**SameResult** This approach groups adjacent states together, such that, whenever a Q-update occurs for a certain Qbert position, all possible states that could have directly led to that state also get updated. The rationale here is that, for most states in the Qbert world, the path Qbert takes to get to that state is usually not important.

Many different state representations were used with varying results. First, a state representation was used keeping track of the color of every block, the position of disks, Qbert’s position, the position of enemies (purple agents) and the position of friendlies (green agents). While this state representation completely encapsulates Qbert’s world, it is very verbose and intractable if not using any distance metric, as described previously.

### 2 Results of generalization

The results of using the previous three distance metrics can be seen in...

### 3 Description of approach to exploration

First, the  $\epsilon$ -greedy approach to exploration was implemented, where a random action is chosen with

probability  $\epsilon$  and the greedy action is chosen otherwise. An  $\epsilon$  value of 10 % was chosen. This is the approach chosen by many researchers using reinforcement learning in the ALE [1, 2, 3]. This method is attractive for many reasons. First, it is very simple to implement with any standard random number generator. Second, despite its simplicity, it can lead to very powerful results, since it can allow an agent to explore states it would never normally explore. However, convergence to an optimal policy can be slow with  $\epsilon$ -greedy. Also, this method can be dangerous in the Qbert game, where a bad random action can lead to sudden death.

Next, weighting with optimistic priors was implemented, using  $N(s, a)$ , i.e., the number of times action  $a$  has been attempted in state  $s$ . The simple exploration function used can be seen in Equation (1). This method can often lead to meaningful convergence much faster than  $\epsilon$ -greedy, since it will directly explore unvisited states, instead of choosing randomly. This method is very attractive in the Qbert world, since exploring unexplored states usually means going on blocks without the goal color, turning them to the correct color.

$$f(u, n) = \begin{cases} 100 & n < 1 \\ u & n \geq 1 \end{cases} \quad (1)$$

### 4 Results of exploration

The random exploration results can be seen in...  
The weighted prior results can be seen in...

### 5 Agent Performance

To learn individual strategies, the subsumption approach described by Mahadevan and Connell was employed [4]. The agent is separated into three learners: an “enemy avoider”, a “friendly finder” and a “block finder”. The structure is shown in Figure 1, where the “S” suppressor nodes denote priority, i.e., avoiding enemies takes priority over catching friendlies and going on blocks. This approach allows the agent to learn various aspects of the Qbert game independently, with a reduced state space for each task. This separation also makes intuitive sense, since learning how to avoid enemies should have little to do with finding blocks.

The enemy avoider is active whenever an enemy (“Coily” and “Purple Ball”) is close to Qbert, i.e. in any of the blocks directly adjacent to Qbert or 2 blocks away. An important strategy is to kill the purple snake enemy (“Coily”) by jumping on to the floating discs on the left and right of the board. If Coily is nearby, he will jump off the board, giving

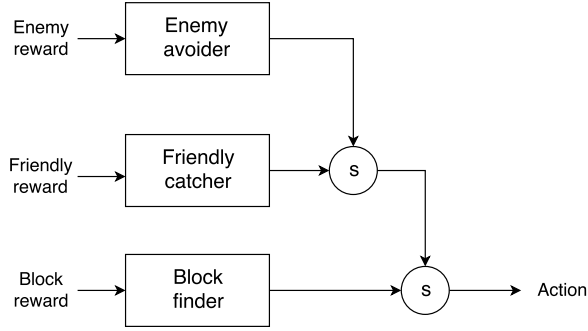


Figure 1: Subsumption network used by the agent.

Qbert 500 points. This score is fed as a reward to the enemy avoider learner. As a penalty for dying to Coily, a negative reward of  $-100$  is given to the learner. This combination of positive and negative reinforcement allows the agent to learn how to avoid and kill Coily.

The friendly catcher is active whenever a friendly green agent (“Sam” and “Green Ball”) is within 2 blocks of Qbert. This learner encourages Qbert to catch Sam for 300 points and Green Ball for 100 points. Stopping Sam is also important because he changes the colors of blocks back to their original values, making Qbert do more work.

The block finder is always active and keeps track of the colors of blocks within 2 blocks’ distance of Qbert. It specifically keeps track of the number of adjacent blocks that are not of the desired color, encouraging Qbert to change them to their correct color to eventually complete the level.

Note that, although one learner’s decision always has priority for executing an action, multiple learners can still learn simultaneously. For example, while Qbert is avoiding an enemy, the block finder learner can still be learning from the actions being taken.

## Acknowledgments

There was a discussion with Andrei Purcarus and Andrew Lowther concerning the positions of various important bytes in the ALE RAM, including the bytes indicating a safe move update and the byte indicating a level change.

## References

- [1] Aaron Defazio and Thore Graepel. A comparison of learning algorithms on the arcade learning environment. *arXiv preprint arXiv:1410.8620*, 2014.

- [2] Marc G Bellemare, Yavar Naddaf, Joel Veness, and Michael Bowling. The arcade learning environment: An evaluation platform for general agents. *J. Artif. Intell. Res.(JAIR)*, 47:253–279, 2013.
- [3] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.
- [4] Sridhar Mahadevan and Jonathan Connell. Automatic programming of behavior-based robots using reinforcement learning. *Artificial intelligence*, 55(2-3):311–365, 1992.