

---

# Autonomous Agents

## Report Assignment 2: Single Agent Learning

---

*Authors:*

Agnes VAN BELLE (10363130),  
Maaïke FLEUREN (10350470),  
Norbert HEIJNE (10357769),  
Lydia MENNES (10333843)

October 3, 2012

## 1 Introduction

This report has been written for the Master Artificial Intelligence course Autonomous Agents. This report will contain the answers, motivation and explanation for our implementations of the tasks we had to accomplish in our second assignment for this course. These tasks were centered around the topic of ‘Single Agent Learning’.

### 1.1 The environment

In all tasks there is assumed to be a grid world (of  $11 \times 11$ ) with a predator and a prey in it. The agents can both move one tile forward each iteration. The direction they take (or if they move at all) is affected by probabilities (their policies). If they move over the edge of the grid they end up at the opposing side of the grid. The prey will never step onto the predator. We are focused on improving the decisions of one agent, the predator.

### 1.2 The state space representation

In the experiments described the first report, we used a state space that was an intuitive, yet cumbersome representation. We will refer to this state space representation as the ‘default’ state space. The amount of states that was used in this default state space was  $(11 \times 11) \times (11 \times 11) = 121 \times 121 = 14641$ .

As can be seen in Figure ??, there is a symmetry in the state space, and thus relatively much values are redundantly computed. We elaborated on a more efficient representation for the states, and finally contrived one consisting of 21 states, an approximately 697 times smaller one than the default state space. We will refer to this state space representation as the “efficient” state space. By using the symmetry that was present in the state space a much smaller state space was achieved. The expectation is that the computational cost will be lowered, since only 21 state values have to be computed instead of 14641 states. The algorithm has not changed however, so it is also expected that we do not see a large difference in the amount of iterations using the same parameters.

This efficient state space has been implemented for both Value Iteration and Policy Evaluation.

Each state represents a distance between the prey and predator. These are represented in the lower left diagonal of a matrix, in which the x-axis is the relative horizontal distance in the MDP and the y-axis the relative vertical distance in the MDP. This matrix is shown in Figure ?. Combinations of positions of prey

and predator for which the horizontal and vertical distances are equal are now treated equivalent. Also two combinations for which the horizontal distance in one equals the vertical distance in the other and vice versa are considered equal. In order to navigate through this state space different actions are required. These are: *approach horizontal*, *retreat horizontal*, *approach vertical*, *retreat vertical* and as before *wait*. When interacting with the environment these actions are converted into corresponding actions in the real world. This only requires the relative direction of the prey (which is always located at the centre, regardless of its coordinates) with respect to the predator. This is computed by using the difference in location of the prey and predator on the x- and y-axis.

This reduction in the number of states has lead to output values which can be viewed in the appendix section ???. The results show that the algorithm does not converge faster in terms of iterations with the efficient state space than with the default state space with  $\theta = 0$  and  $\gamma = 0.8$ , the default state space required 111 iterations and the efficient state space required 107.

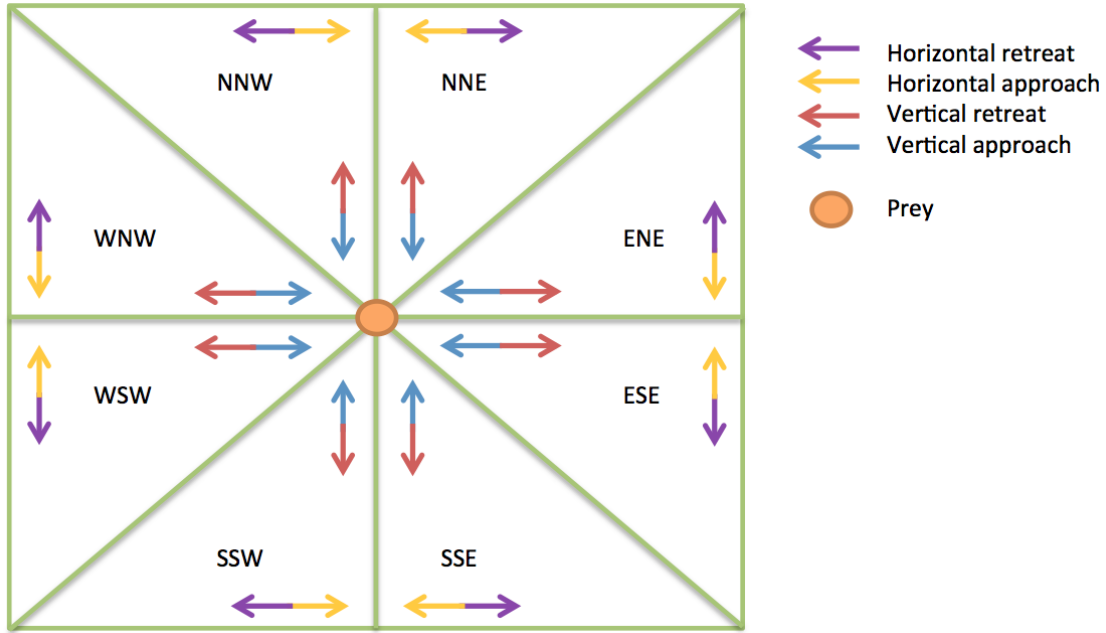


Figure 1: beschrijving

### 1.3 Implementation details

This report will not be about our exact code and implementation details. However, a class diagram of our code is provided in Appendix A.

## **2 Learning algorithms**

### **2.1 (M) Q-Learning**

### **2.2 (M) Experiment with $\epsilon$ and optimistic initialization**

### **2.3 (SC) Softmax action selection instead of $\epsilon$ -greedy**

### **2.4 (SC) On-policy Monte Carlo Control**

### **2.5 (SC) Off-Policy Monte Carlo Control**

### **2.6 (SC) Sarsa**

## **3 Conclusion**

# Appendices

## A Class Diagram