

Assignments Autonomous Agents

D.M. Roijers Deepak Viswanathan

October 13, 2014

1 Rules and schedule

In this course, 50% of your grade will be determined by the practical assignments. There are three practical assignments, each of which weighs 33% for the lab part of the course. The assignments are *Single Agent Planning*, *Single Agent Learning*, and *Multi-Agent Planning and Learning*. All of these assignments are provided in this document.

The assignments should be done in groups of three or four. Each assignment consists of programming and writing a report. The code and the report should be handed in together (e.g. in a zip file), at the given deadlines.

1.1 Quality demands

The report for each assignment must be written in correct academic English. Points can and will be deducted if this requirement is not met. All reports must have an introduction and a conclusion. References must be provided whenever used. We assume all students are familiar with proper citation. If this is not familiar to you, please contact the TA; plagiarism will be punished.

The code should be readable and commented. It should be clear from the comments, which function provides which functionality. If you are using an object oriented language, be sure to comment on the class structure as well. The names of variables and classes should be meaningful; do not use 1, 2 or 3 character variable names except for in counters. When in doubt about the amount of comments you should add to your code, use as a guideline that the comments should lead the reader through your code.

Any programming language may be used for programming the assignments. However, the usage of esoteric and low-level languages (like assembly) is strongly discouraged, as this hardly ever leads to readable code (which is required).

1.2 Schedule

- Assignment 1 is due November 14, 2014, at 11:59PM.
- Assignment 2 is due November 28, 2014, at 11:59PM.
- Assignment 3 is due December 12, 2014, at 11:59PM.

Turning in assignments late is not allowed. Assignments that are turned in late will not be graded. The time your work arrives at the inbox of the TA counts. We recommend submitting a couple of hours before the deadline, to avoid last minute system failure issues.

1.3 Grading

Every assignment will consist of different functionalities you have to implement and test. There are two types of functionalities. The *must-have* (M) functionalities *must* be correctly implemented - they form the basic functionality of your program, and the baseline for comparisons against other methods. It is impossible to pass an assignment without correctly having implemented the *must-haves*. Then there are *should/could-haves* (SC), which are functionalities that are less essential to the assignments. Implementing SCs will get you extra points for your grade. Implementing only the Ms will get you a 5.5 out of 10, iff done correctly. Implementing enough SCs, after implementing all the Ms, can increase your grade to a 10 out of 10. We strongly recommend implementing all the Ms first. Having all SCs, but not all the Ms will automatically give you a bad mark (failing the assignment), because you would miss the essential functionality, and the base-lines to compare against.

2 Assignment 1: Single Agent Planning

2.1 Intro

In this assignment we are going to use a predator vs prey grid world MDP. In this assignment there will be one predator, and one prey. In later assignments, we will add more predators.

The basic environment is a 11 by 11 grid (see Figure 1). The grid is toroidal, which means that the north side of the grid is attached to the south side, and the east side is attached to the west side. So the adjacent squares to the square with coordinates $(0,0)$, are $(1,0)$, $(0,1)$, $(10,0)$ and $(0,10)$. The predator and the prey can be anywhere on this grid, though they cannot be on the same square. The starting position of the prey is $(5,5)$ and that of the predator $(0,0)$.

The prey behaves in a known way, and can therefore be modelled as part of the environment. The prey stays at the same square with probability 0.8, and moves randomly, and with equal probability to any of the adjacent squares (North, West, South or East from the current position) that are open, at every time-step. Therefore, if time starts at $t = 0$, at time-step $t = 1$, the prey has a 0.05 probability of being on $(6,5)$, $(5,6)$, $(4,5)$, and $(5,4)$, and a 0.8 probability of being at $(5,5)$. Note that the prey would never move into the predator, so the probabilities are different when the prey is standing next to the predator.

The predator is the agent. After implementing the environment, you will mainly focus on writing code for the predator agent. The predator has five possible actions: *North* which results in the predator being moved upwards, a difference of $(-1, 0)$ in coordinates. *South*, moving the predator downwards: a $(+1,0)$ change in coordinates. *East*, $(0,+1)$, and *West*, $(0,-1)$. The final possible action is *Wait*, for which the predator does not move.¹ The goal for the predator is to catch the prey. This happens when the predator stands next to the prey, and moves in the direction of the prey. (The prey has no chance of escape in this case.) When the prey is captured, the episode ends, and the game is reverted to the starting position.

The immediate reward for catching the prey is 10, the immediate reward for anything else is 0.

¹Note that this might not seem like a very useful action, but it can be useful when we move to the multi-agent setting.

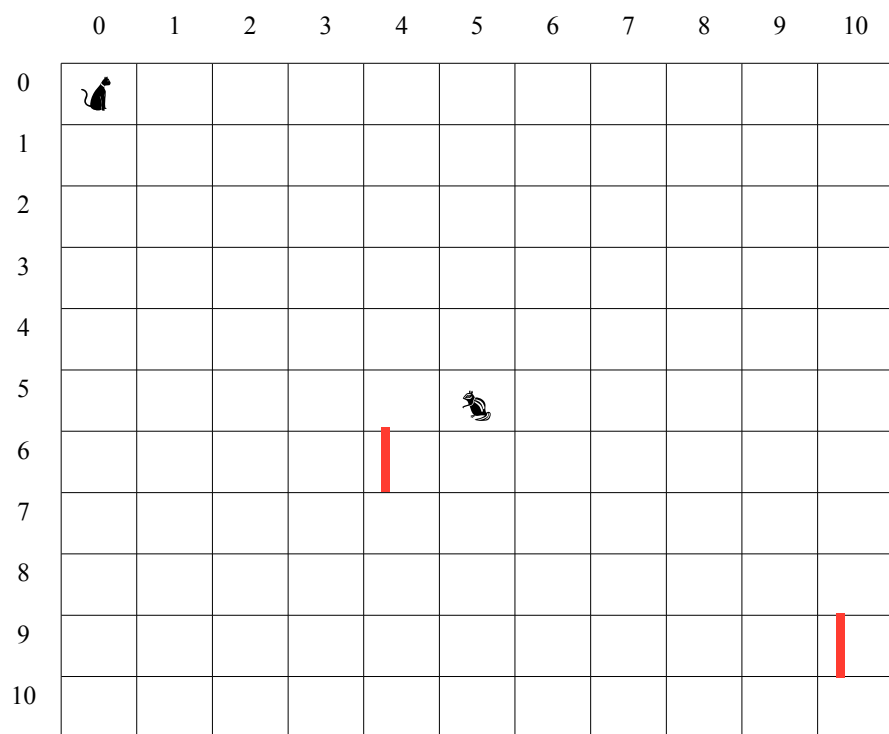


Figure 1: The predator-prey environment: a 11x11 toroidal grid, in the starting position. The predator is depicted by the cat, and the prey by the squirrel.

2.2 Assignments

In Assignment 1, we will consider the planning scenario: the entire MDP is known to your agent. The agent can thus determine the optimal policy even before interacting with the environment. If you are asked to run simulations you should first running your planning algorithm, and then execute the policy that results from it.

The following list contains the sub-assignments for assignment 1. A mandatory assignment is denoted with an **M**. Should/Could-haves are indicated by **SC** followed by a number of points. The points you will get for completing all the must-haves is 55, the other 45 points you can get by completing the Should/Could-haves.

We highly recommend that you show the work you did for the Must-haves (code and results) to your TA, before proceeding to the other assignments.

1. **M** Write a simulator for the environment. Be sure to keep the policies of the predator and the prey separate; both will change in the upcoming assignments. For now, use a random policy for the predator (picking one of the five actions randomly, with equal probability). Test your environment by running a few simulations. Print the states and the actions of the predator, so that you can check whether your simulator is functioning properly. Encode the state as `Predator(X,Y)`, `Prey(X,Y)` (where X and Y are the coordinates). Measure the time it takes on average (for 100 runs) for the predator to catch the prey with this random policy, mention the average and the standard deviation in your report.
2. **SC (10)** Code and use iterative policy evaluation² to determine the value of the random policy for all possible states, using a discount factor of 0.8. Put the values of the following states in your report:
 - `Predator(0,0)`, `Prey(5,5)`
 - `Predator(2,3)`, `Prey(5,4)`
 - `Predator(2,10)`, `Prey(10,0)`
 - `Predator(10,10)`, `Prey(0,0)`

Also report how many iterations it takes to converge.

3. **SC (10)** Implement Policy Iteration³, and output the values of all states in which the prey is located at (5,5). Test the convergence

²see <http://webdocs.cs.ualberta.ca/~sutton/book/ebook/node41.html>

³See <http://webdocs.cs.ualberta.ca/~sutton/book/ebook/node43.html>

speed (in number of iterations), for different discount factors: at least 0.1, 0.5, 0.7 and 0.9, and report on the results. Compare the results to those of Value Iteration.

4. **M** Implement Value Iteration⁴, and output the values of all states in which the prey is located at (5,5). Test the convergence speed (in number of iterations), for different discount factors: at least 0.1, 0.5, 0.7 and 0.9, and report on the results.
5. **SC (25)** Think of a smarter way to encode the state-space such that the size of the state-space is reduced as much as possible. How much does your solution change the state-space? How does this effect the runtime of the algorithms?

NB: this assignment is highly recommended, as it will lead to great speed-ups in the following assignment as well.

3 Assignment 2: Single Agent Learning

In Assignment 2, we will use the same environment as in Assignment 1, but now we will assume the learning scenario: the agent does not know the transition probabilities, nor the reward structure. On a very high level there are two ways to come to a good solution in this setting: learning the model, and do planning again (model based learning), or not learn the model, and directly try to learn a high-reward policy (model-free learning). In this assignment we will focus on the latter.

1. **M** Implement Q-learning⁵ and use this for your predator agent. Use ϵ -greedy action selection with $\epsilon = 0.1$. Initiate the values for your Q-learning table optimistically with a value 15 for all cells in the table. Show plots on the performance of the agent over time for different α (at least 0.1, ..., 0.5), for different discount factors (at least 0.1, 0.5, 0.7 and 0.9).
2. **M** Experiment with different values of ϵ and the optimistic initialization of the Q-table. Make up good values to test, and explain why you chose these values.

⁴See <http://webdocs.cs.ualberta.ca/~sutton/book/ebook/node44.html>

⁵See <http://webdocs.cs.ualberta.ca/~sutton/book/ebook/node65.html>

3. **SC (10)** Use Softmax action selection⁶, instead of ϵ -greedy. And illustrate the difference, using graphs from your empirical results.
4. **SC (15 points per method)** Implement other ways to do learning, instead of Q-learning, from the following list, explain the difference theoretically, and compare the different methods using informative graphs. Note that your grade cannot become higher than 10 out of 10.
 - On-policy Monte-Carlo Control⁷
 - Off-Policy Monte Carlo Control⁸
 - Sarsa⁹
 - If you would like to test a different learning method, please discuss this with your TA. The TA will have to approve your proposed method. (We highly recommend not to start implementing this until you have received this approval.)

4 Assignment 3: Multi-Agent Planning and learning

In this assignment we will make several changes to the environment and the agents. We will make the prey intelligent, making it harder to catch.

Firstly, both the prey and the predator will be agents. That means that you will apply your algorithms to both the prey and the predator.

Secondly, we change the way the agents move on the grid. The prey and the predator will take an action simultaneously. If the result of both actions is that the prey and the predator end up on the same square, the predator gets a +10 reward, and the prey gets a -10 reward. Every other transition results in a reward of 0 for both predator and prey.¹⁰

Thirdly, in some of the optional assignments, we can add extra predators to the grid at (10, 10), (10, 0) and (0, 10). (This might look like the predators are far apart, but actually they are back-to-back.) If two predators run into each other – end up at the same square after their actions – this creates such a confusion, that the prey escapes. The predators get a reward of -10 and

⁶See <http://webdocs.cs.ualberta.ca/~sutton/book/ebook/node17.html>

⁷<http://webdocs.cs.ualberta.ca/~sutton/book/ebook/node54.html>

⁸<http://webdocs.cs.ualberta.ca/~sutton/book/ebook/node56.html>

⁹<http://webdocs.cs.ualberta.ca/~sutton/book/ebook/node64.html>

¹⁰Please note that the requirement that the prey never moves into the predator from the previous assignments no longer exists.

the prey a reward of +10. These too are absorbing states, and the episode ends.

Finally, because the prey is now an agent, it can actively take a Wait, North, East, South or West action. The results of which are almost the same as for the predators (see previous assignment), with one important difference: there is a probability the prey trips, in which case the prey stays where it is, regardless of the action it took. The probability of tripping is $p_{trip} = 0.2$.¹¹

Note that the predators are a team. They all obtain the same reward. If two predators end up on the same square, the reward is -10 and the episode ends, even if one of the predators catches the prey in that timestep as well. We only look at where the agents end up after they have taken their actions. So if two predators “pass eachother” (e.g. one moves from $(0, 0)$ to $(0, 1)$ and the other moves from $(0, 1)$ to $(0, 0)$), the episode does not end.

4.1 Assignments

1. **M** Implement the new environment, adding an option to include 1, 2, or 3 more predators. Run a simulation with random policies to check whether the environment works correctly.
2. **M** Choose one of the following two options:
 - Apply independent Q-learning¹² and analyze what happens for 1, 2, 3 and 4 predators. Note that the prey learns also. Illustrate the performance of the this algorithm for different parameter settings, with informative graphs.
 - **(+10 points)** Read the paper “Markov games as a framework for multi-agent reinforcement learning” by M. L. Littman (1994)¹³, summarize this article briefly, implement the minimax-Q algorithm (Figure 1 in the paper), and apply it to the 1 predator, 1 prey scenario. Analyze the performance of the algorithm, for different parameter settings. Use informative graphs to illustrate the performance of this algorithm in your report.

¹¹Note that the tripping probability is necessary, especially when there is only one predator. If this probability would not exist, the prey could run away succesfully indefinitely.

¹²“A Concise Introduction to Multiagent Systems and Distributed Artificial Intelligence” by Nikos Vlassis (2007), Section 7.3.1

¹³This paper is available at <http://www.cs.rutgers.edu/~mlittman/papers/refer.html>

Doing both of the options, and comparing the results will give you 85 points instead of 55 for the must-haves.

3. **SC (45)** Implement other learning (or planning) algorithms for Markov Games, and analyze the performance of these algorithms in your report. You can freely use the different possibilities for the environment. Changes to the environment are also allowed, as long as you explain them well. Make sure to properly cite the sources you are using in your report. Points will be awarded based on an assessment of the difficulty of the chosen method, as well as the quality of the analysis. If you want to know how many points are maximally possible with the method of your choice before implementing it, please contact the TA.

5 Contact info

Please contact your TA for any additional questions.

Email: D.GeethaViswanathan@uva.nl.