

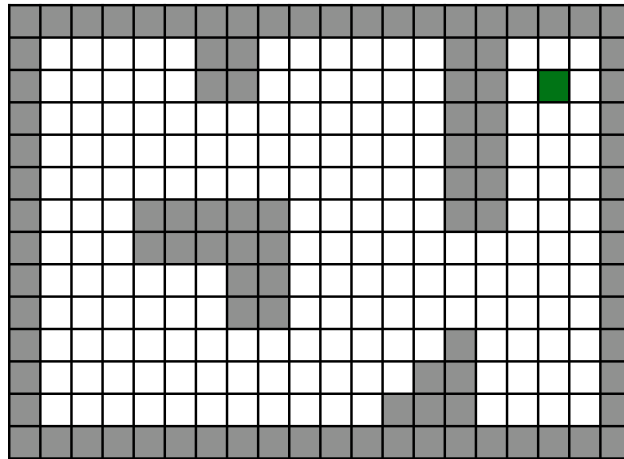
# Autonomous Robots

## Robot Learning – Reinforcement Learning

This document will guide you through the practical work related with the Reinforcement Learning subject of the Autonomous Robots course. The goal of this practical exercise is to implement a Reinforcement Learning algorithm to learn a policy that moves a robot to a goal position. The algorithm is the Q-learning algorithm. You will program all the code in Matlab.

### 1. The reinforcement learning problem

The problem consists in finding the goal in a finite 2D environment that is closed and contains some obstacles (in grey) as shown in the figure:



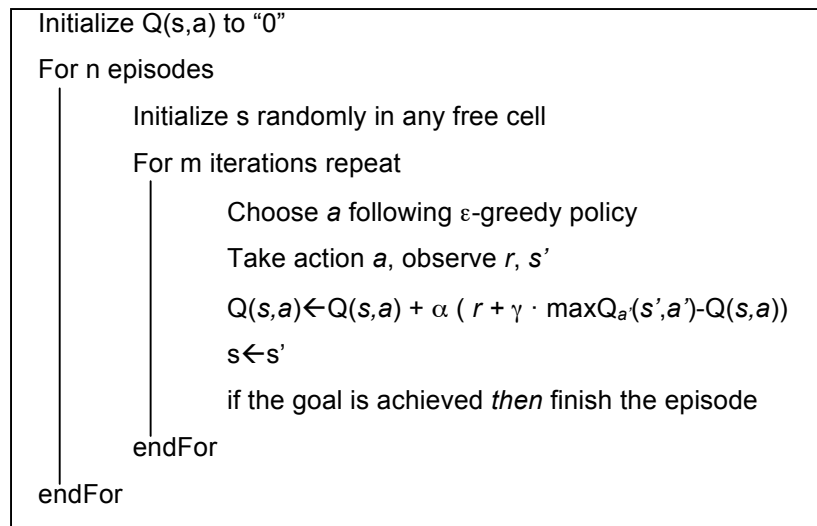
**States and actions:** The size of the environment is  $20 \times 14 = 280$  states. The robot can only do 4 different actions:  $\{\leftarrow, \uparrow, \rightarrow, \downarrow\}$  (not diagonal movements!). Therefore the size of the Q function will be  $280 \times 4 = 1120$  cells.

**Dynamics:** The robot can be located in any free cell (not in the obstacle cells!). The function that describes the dynamics is very simple: the robot will move ONE cell per iteration to the direction of the action that we select, unless there is an obstacle or the wall in front of it, in which case it will stay in the same position.

**Reinforcement function:** Since the goal is to reach the goal position as fast as possible, the reinforcement function will give -1 in all cells except in the goal cell, where the reward will be +1. The cell that contains the goal is (18,3).

## 2. The algorithm

You must implement the Q-learning algorithm. Consider the problem as episodic and repeat the learning process during hundreds of episodes. The algorithm should be like:



You will have to set several parameters experimentally: n, m,  $\epsilon$ ,  $\alpha$  and  $\gamma$ .

## 3. Testing the learnt policy

In order to test the learnt policy, after some episodes you can test the policy. Just run an episode with  $\epsilon = \alpha = 0$ , and sum the rewards. Since the initial position is randomly selected, run several episodes and compute the average of the accumulated rewards. This average should decrease with respect to the number of episodes in which the learning was enabled.

You can also plot the greedy policy by checking for each of the 280 cells the best action. Then you can determine if the best action is the one you would choose.

## 4. Implementation

You must program a Matlab function that has the following definition. It is very important that the input and output arguments have the following format.

```
Q=q_learning(map,q_goal,alpha,gamma,epsilon,n_episodes,n_iterations)
```

where:

map: matrix that represents an occupancy grid map

q\_goal: X and Y coordinates of the cell that contains the goal state.

alpha: learning rate of the algorithm

gamma: discount factor of the algorithm

epsilon: random action probability

n\_episodes: number of episode repetitions

n\_iterations: number of iterations per episode

Q: 3D matrix of size: number of columns; number of rows; number of actions. The matrix contains the Action Value Function.

When your function is called with the following arguments,

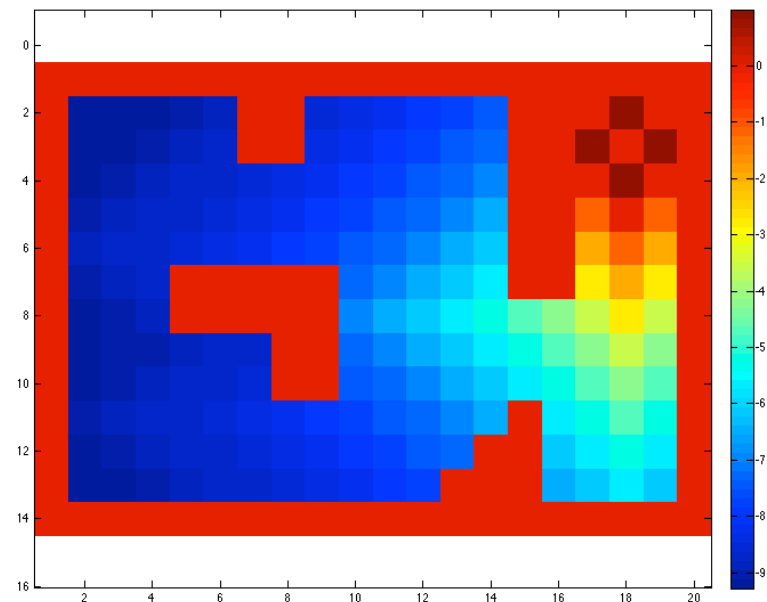
```
Q=q_learning(map,[18,3],0.1,0.9,0.3,20000,50)
```

a) Q function:

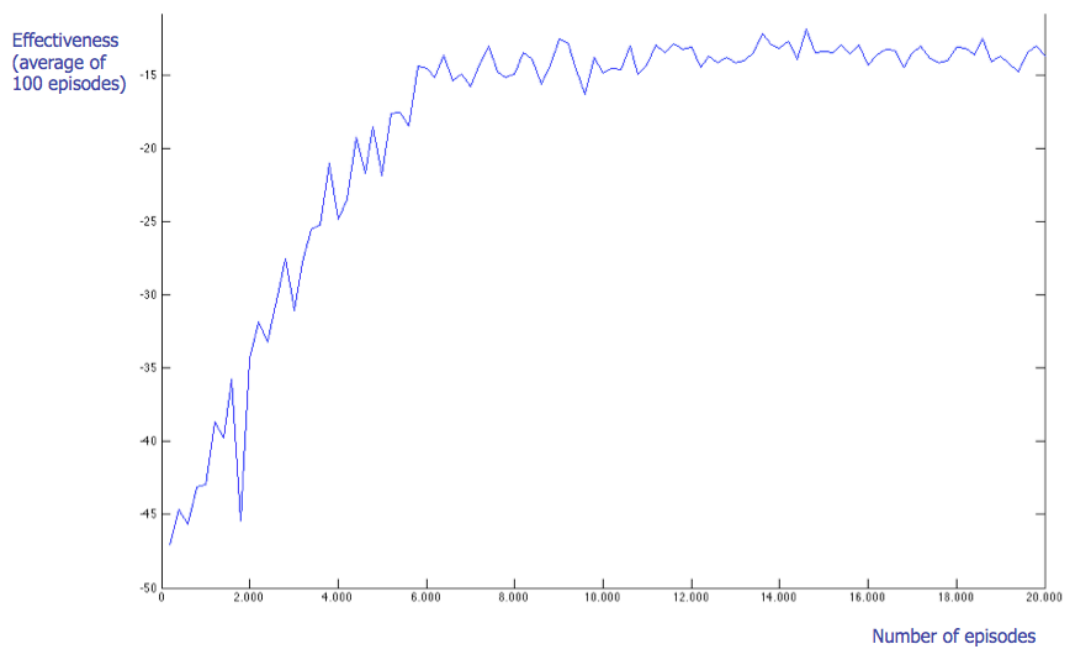
b) The optimal policy using some kind of representation such as:

o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o
o	v	>	>	>	v	o	o	>	>	>	>	>	v	o	o	v	^	v	o
o	>	>	>	>	v	o	o	>	>	>	>	>	v	o	o	>	g	<	o
o	>	>	>	>	>	>	>	>	>	>	>	>	v	o	o	>	^	<	o
o	>	>	>	>	>	>	>	>	>	>	>	>	v	o	o	>	^	<	o
o	>	>	>	>	>	>	>	>	>	>	>	>	v	o	o	>	^	<	o
o	^	^	^	o	o	o	o	o	>	>	>	>	v	o	o	^	^	<	o
o	^	^	^	o	o	o	o	o	>	>	>	>	>	>	>	^	^	^	o
o	v	>	>	>	<	v	o	o	>	>	^	^	^	^	^	^	^	^	o
o	v	>	>	>	>	v	o	o	^	>	>	>	>	^	>	^	^	v	o
o	v	>	>	>	>	>	>	>	^	>	^	^	^	o	>	^	^	^	o
o	>	^	v	>	>	>	v	>	^	^	^	^	o	o	<	^	v	>	o
o	>	<	<	>	>	<	>	>	>	v	o	o	o	o	<	^	v	<	o
o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o

c) A graphical representation of the State Value Function  $V$  such as:



d) A graphical representation of the evolution of the effectiveness such as:



## 5. Submission.

**WORK TO DO:** Submit a report in pdf and the Matlab file “q\_learning.m”. Explain in detail with graphical information, in the report, the work done. Explain also the problems you found.

NOTE that your function “q\_learning.m” will be tested with other environments and other parameters. Remember that your function MUST return the Q function and SHOW GRAPHICALLY the optimal policy, the state value function and the effectiveness evolution.