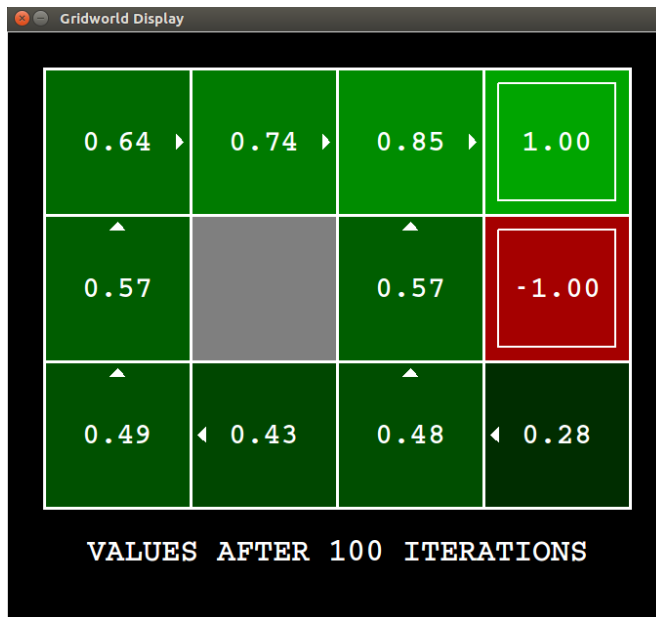ANTALYA BİLİM
ÜNİVERSİTESİ

Assignment 3: **Reinforcement Learning**
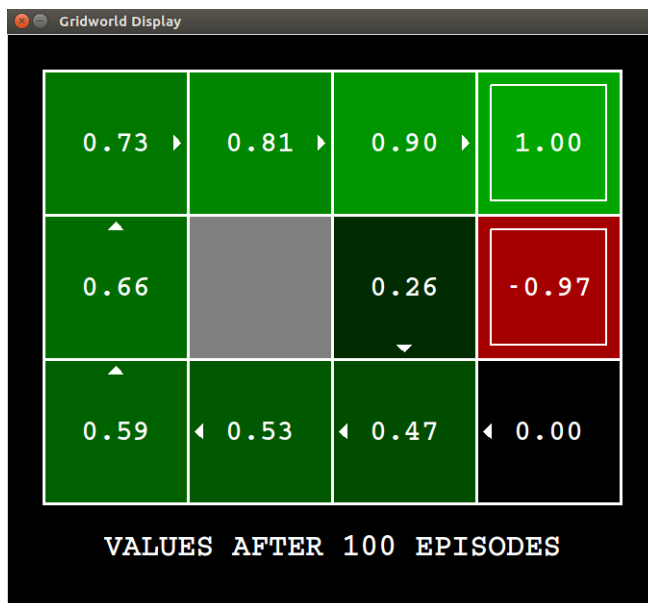
**Ahmed Amine**
**TALEB BAHMED**
**130201120**

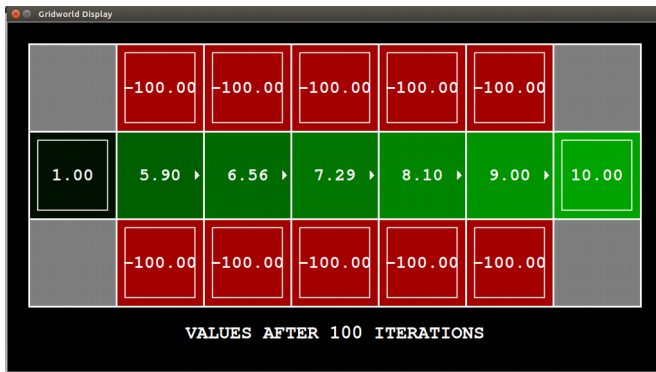**Dr. Hilal KAZAN**

Value Iteration:



1-
The learned values are smaller then those learned by value iteration, since the agent does not have the transition probabilities and the reward for every transition in Q_learning. To get to the optimal value we make the learning rate small enough.



2- Train your Q-learner on the BridgeGrid with no noise (-n 0.0) for 100 episodes. How do the learned q-values compare to those of the value iteration agent? Why will your agent usually not learn the optimal policy?
The Q-learner values are smaller and insufficient to reach the optimal than the value iteration agent values, we need more episodes

VALUES AFTER 100 ITERATIONS



VALUES AFTER 100 EPISODES

3-
Train your Q-learner on the CliffGrid for 100 episodes. Compare the value it learns for the start state with the average returns from the training episodes (printed out automatically). Why are they so different?



6-
When epsilon is -e 0.9 it learns faster than e 0.1, it acts on current policy while, when epsilon is -e 0.1 learns too slow since it acts randomly.



VALUES AFTER 100 EPISODES