

## My First Reinforcement Learning

- **Number of States:** 70 (7 x 10 grid)
- **Number of Actions:** 4 [Up (0), Right (1), Down (2), Left (3)]

The strategy used to balance the learning rate ( $\alpha$ ) and the exploration rate ( $\epsilon$ ) was to keep them constant. The exploration rate could be implemented with a variation in function of the number of episodes, like:

$$\epsilon = e^{-\text{Number of Episodes}}$$

To design the Q-Learning algorithm we used the pseudocode provided in the Temporal-Difference Learning chapter of the book Reinforcement Learning by Richard S. Sutton and Andrew G. Barto. The former is shown below.

### Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

```
Algorithm parameters: step size  $\alpha \in (0, 1]$ , small  $\epsilon > 0$ 
Initialize  $Q(s, a)$ , for all  $s \in \mathcal{S}^+, a \in \mathcal{A}(s)$ , arbitrarily except that  $Q(\text{terminal}, \cdot) = 0$ 
Loop for each episode:
  Initialize  $S$ 
  Loop for each step of episode:
    Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
    Take action  $A$ , observe  $R, S'$ 
     $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$ 
     $S \leftarrow S'$ 
  until  $S$  is terminal
```

To implement the policy algorithm, we used the explanation made in the same book but for  $\epsilon$ -greedy algorithm:

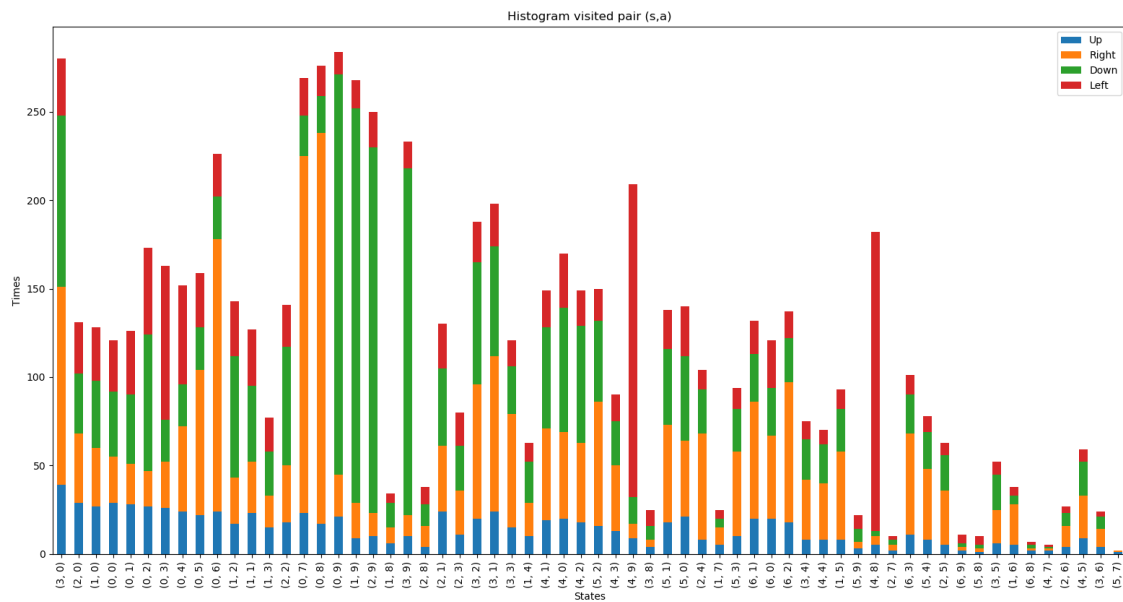
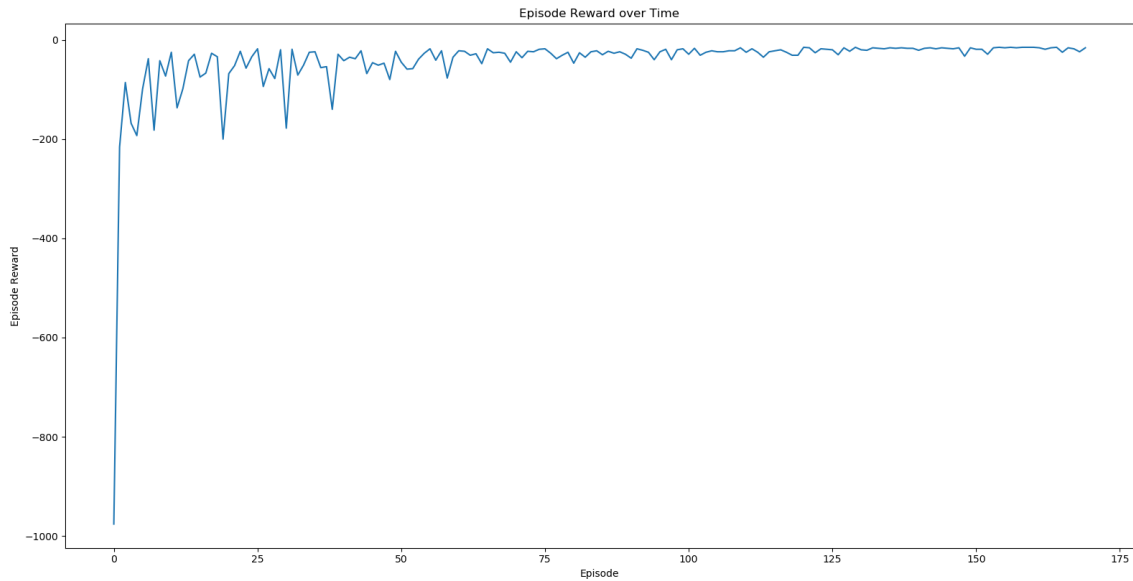
The on-policy method we present in this section uses  $\epsilon$ -greedy policies, meaning that most of the time they choose an action that has maximal estimated action value, but with probability  $\epsilon$  they instead select an action at random. That is, all nongreedy actions are given the minimal probability of selection,  $\frac{\epsilon}{|\mathcal{A}(s)|}$ , and the remaining bulk of the probability,  $1 - \epsilon + \frac{\epsilon}{|\mathcal{A}(s)|}$ , is given to the greedy action.

The states and actions of the environment are given by the Gym library as well as the reward in every state and a “done” flag when the agent steps in the goal state. The Q-Values are stored inside a dictionary in python and are initialized in 0, in other words, every state is initialized with the optimal value, given the current environment in which every state but final, have a negative reward.

The first try made with the implemented algorithm we used the same values for alpha, gamma, and epsilon provided by the solution given in the book:

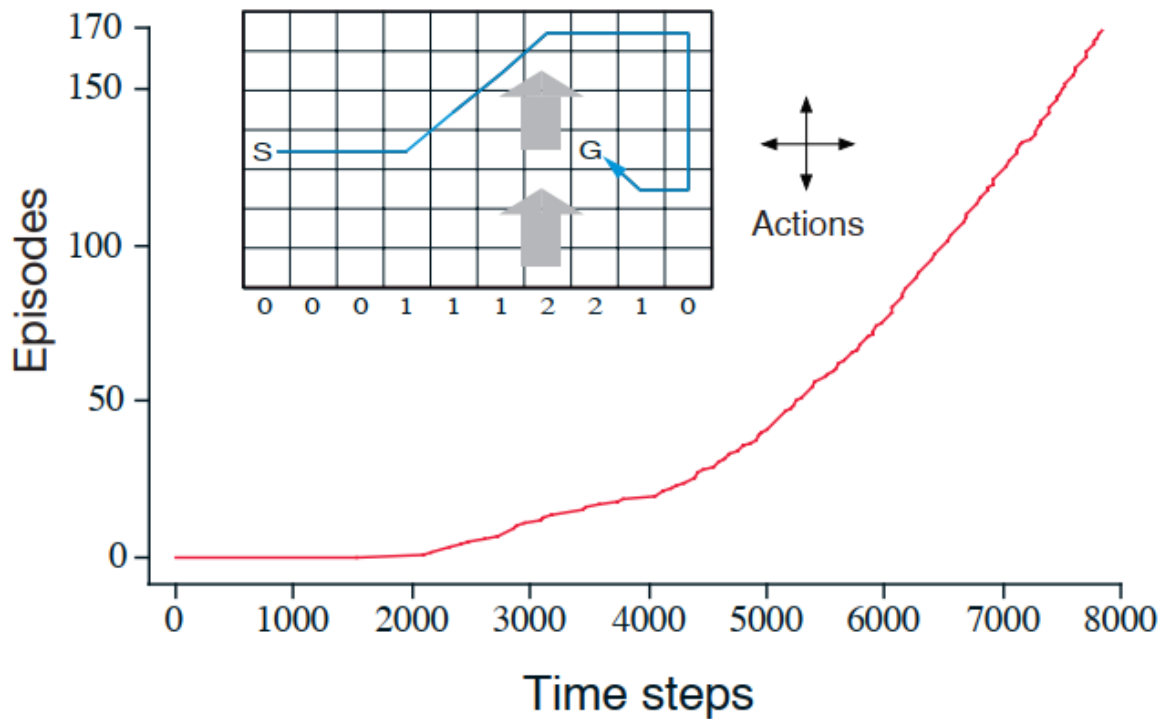
- **Alpha ( $\alpha$ ):** 0.5
- **Gamma ( $\gamma$ ):** 1
- **Epsilon ( $\epsilon$ ):** 0.1
- **Number of Episodes:** 170

The obtained results are shown below:





And as seen, the final result is completely satisfactory, fulfilling the solution given by the book.



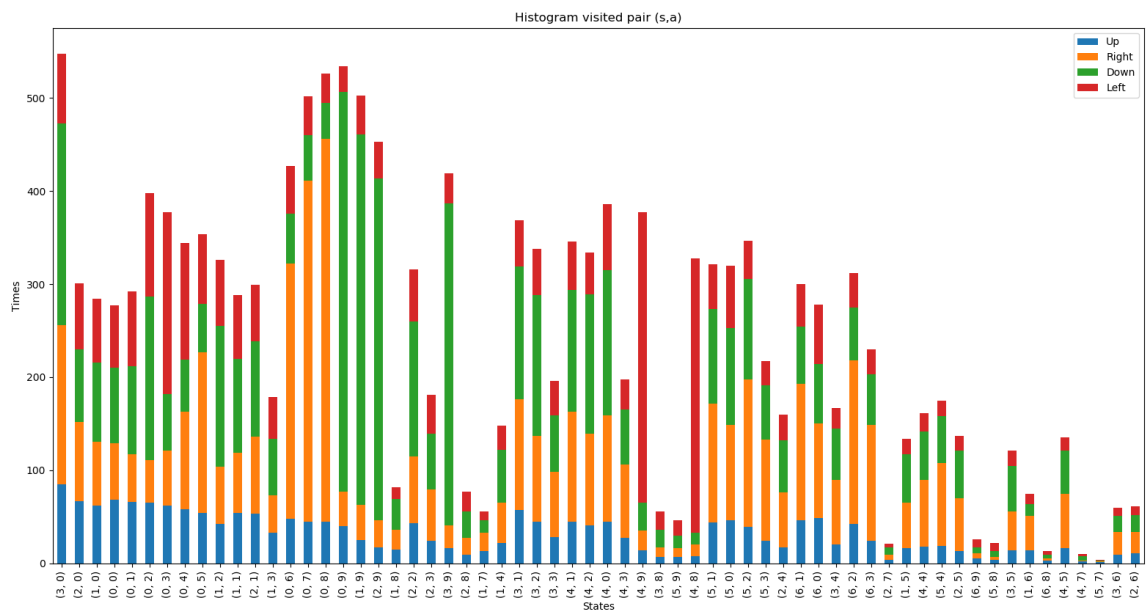
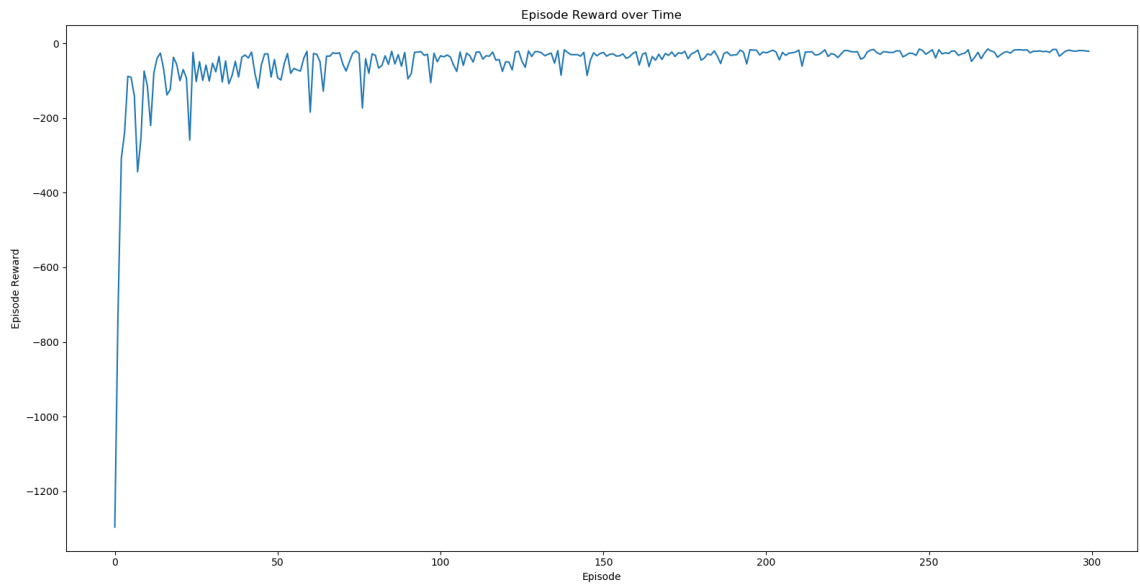
To realize the hyper-parameter sweeping, we decided to sweep 3 times each variable, and in the case of alpha, gamma and epsilon to implement a little piece of code in which the algorithm will use the minimal quantity of episodes to resolve the gridworld, this with the sake of comparison between variations.

**First: Sweeping Alpha (Learning Rate)**

- **Alpha ( $\alpha$ ):** 0.2
- **Gamma ( $\gamma$ ):** 1
- **Epsilon ( $\epsilon$ ):** 0.1
- **Number of Episodes:** 300

**Minimal number of Episodes:** 151, 157, 154, 133, 179

**Average number of Episodes:** 155

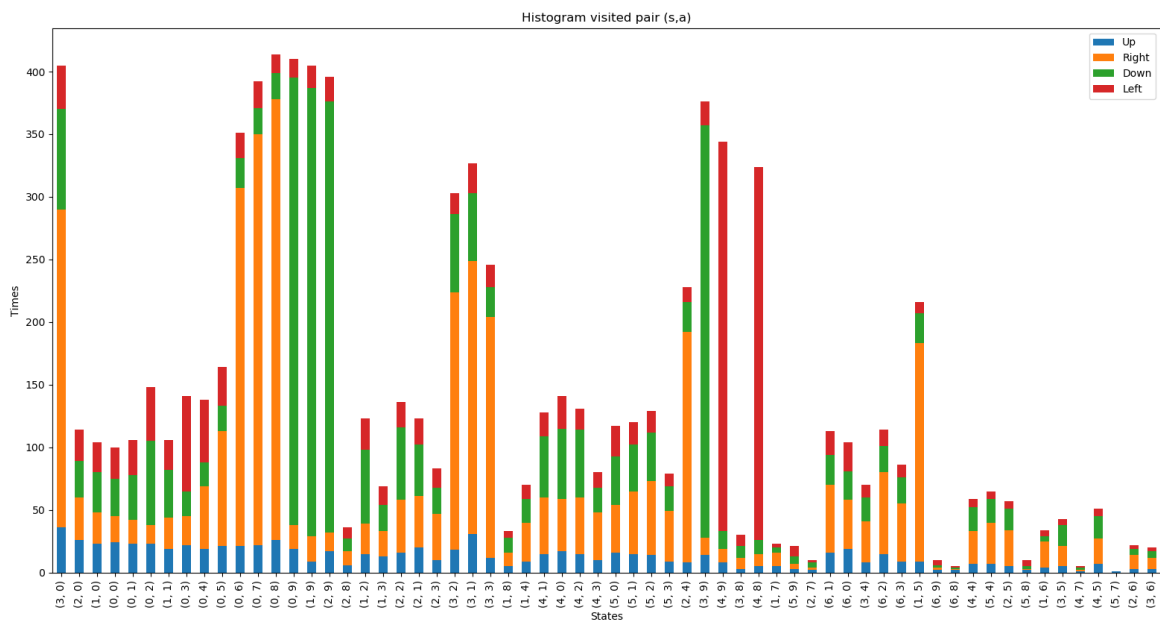
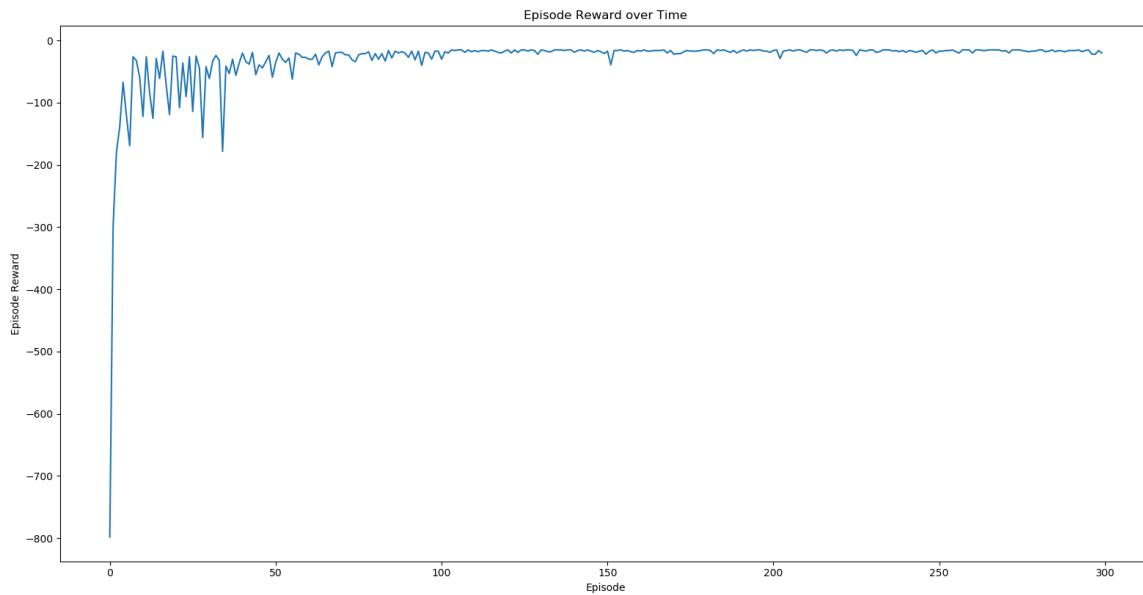


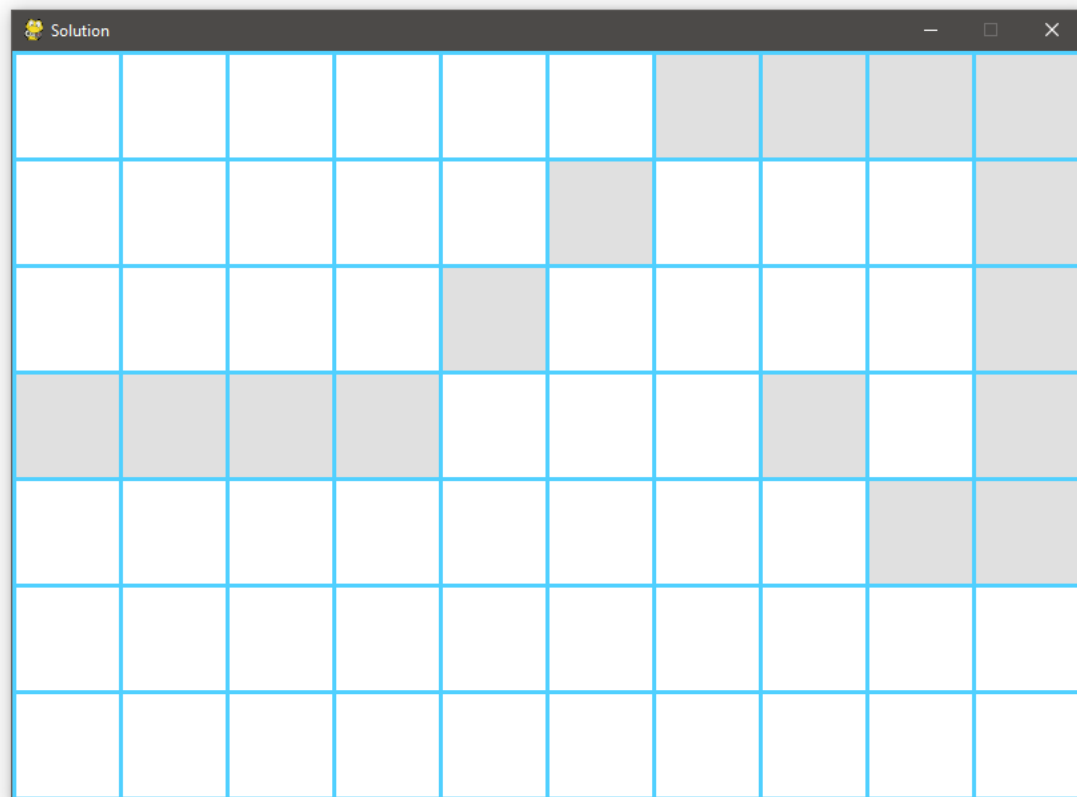
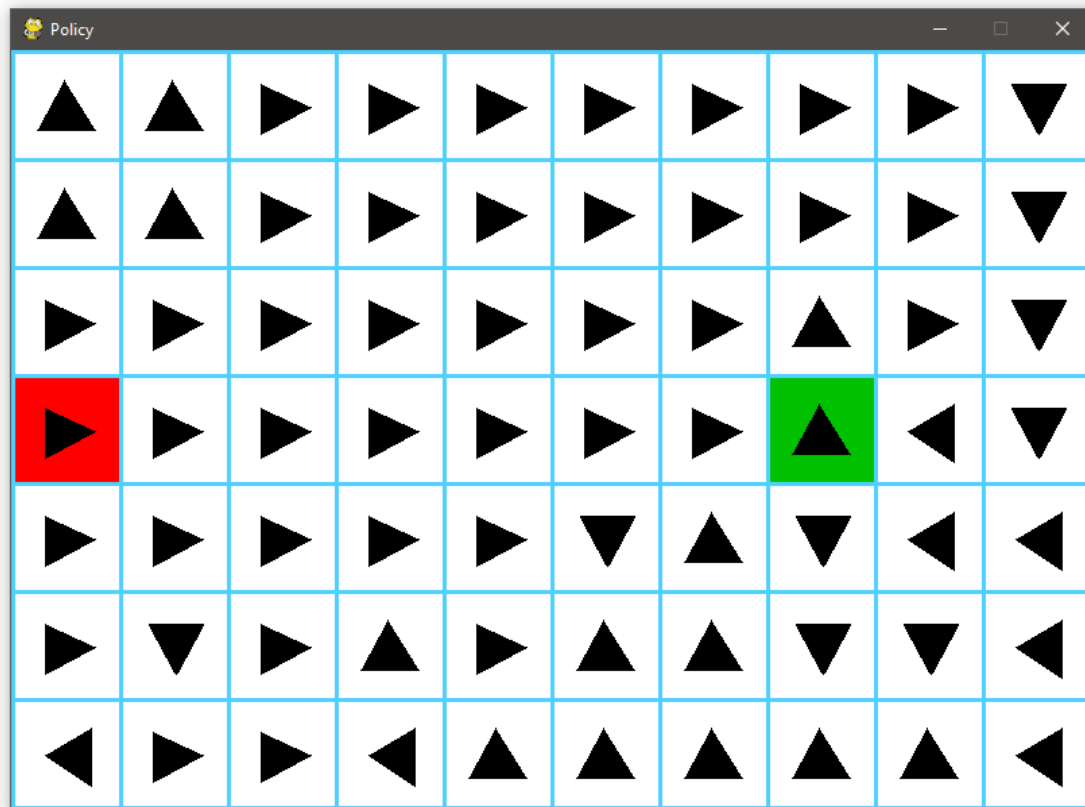


- Alpha ( $\alpha$ ): 0.6
- Gamma ( $\gamma$ ): 1
- Epsilon ( $\epsilon$ ): 0.1
- Number of Episodes: 300

Minimal number of Episodes: 59, 42, 58, 52, 46

Average number of Episodes: 51



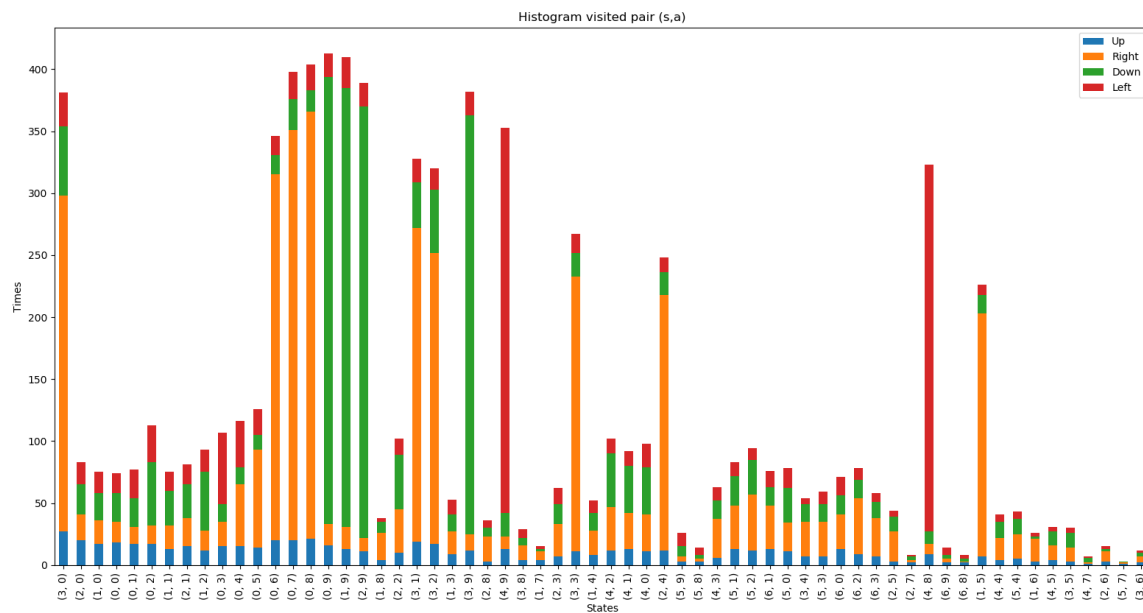
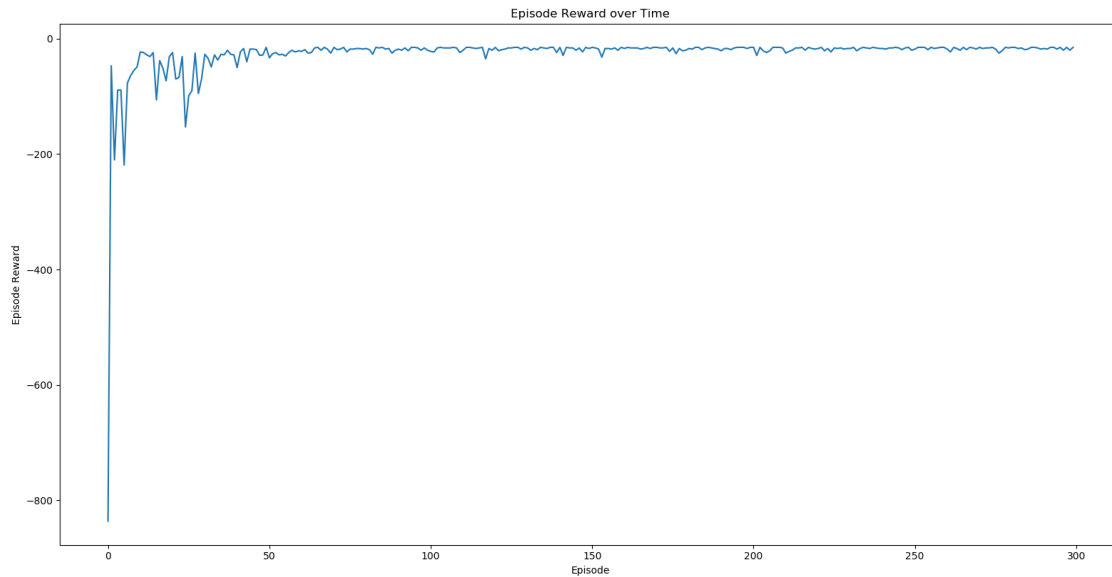




- Alpha ( $\alpha$ ): 0.9
- Gamma ( $\gamma$ ): 1
- Epsilon ( $\epsilon$ ): 0.1
- Number of Episodes: 300

Minimal number of Episodes: 50, 58, 37, 62, 55

Average number of Episodes: 52



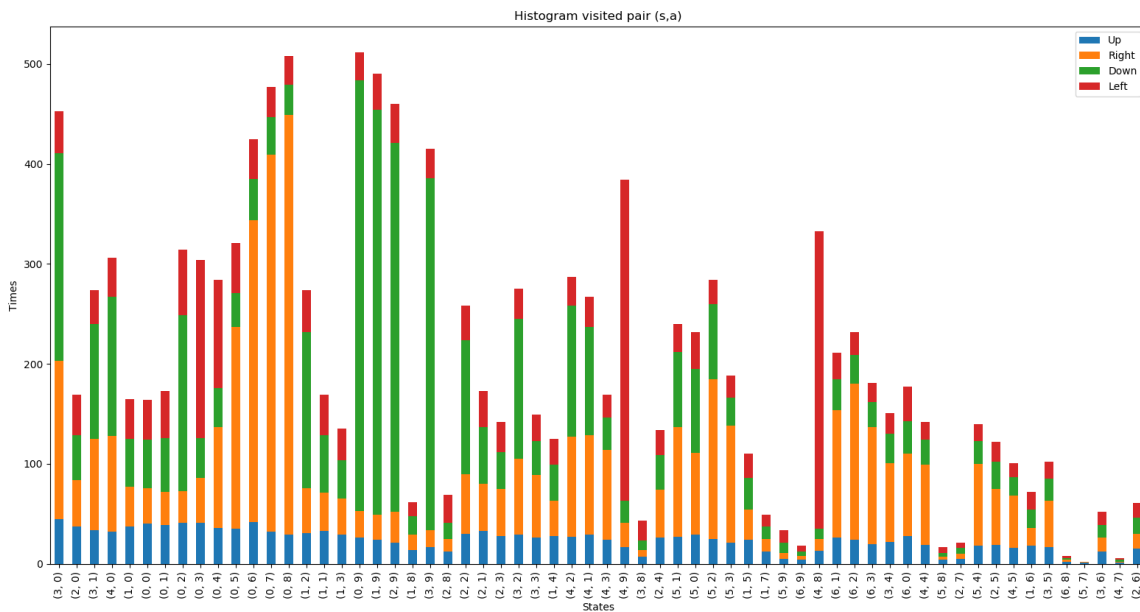
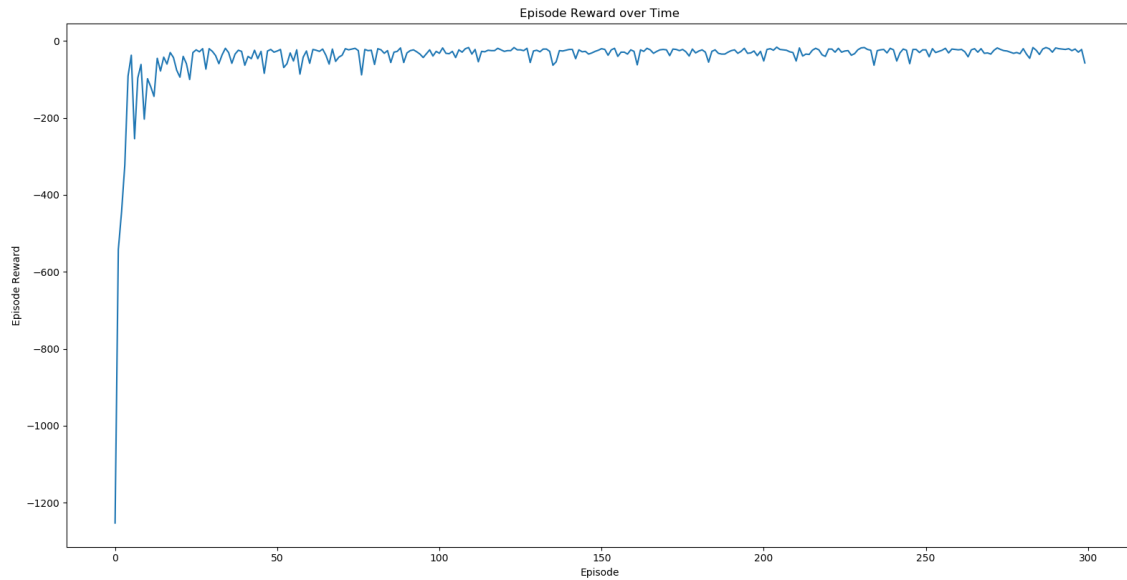


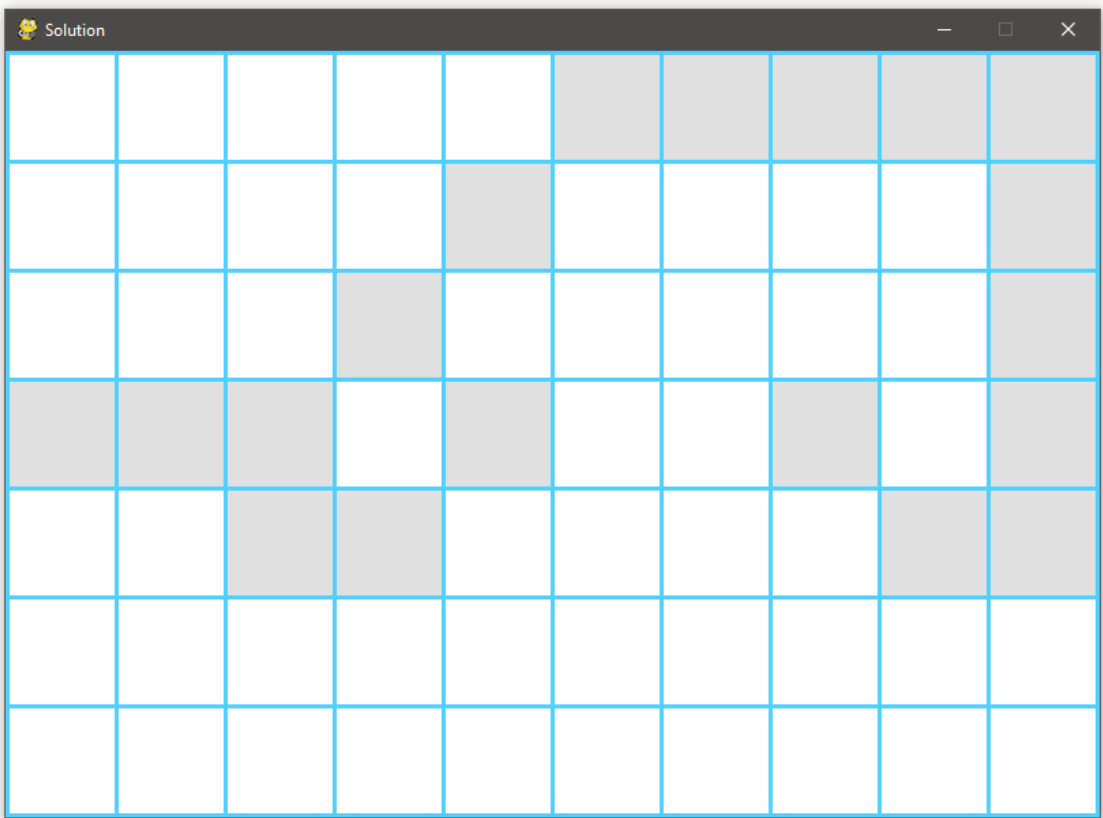
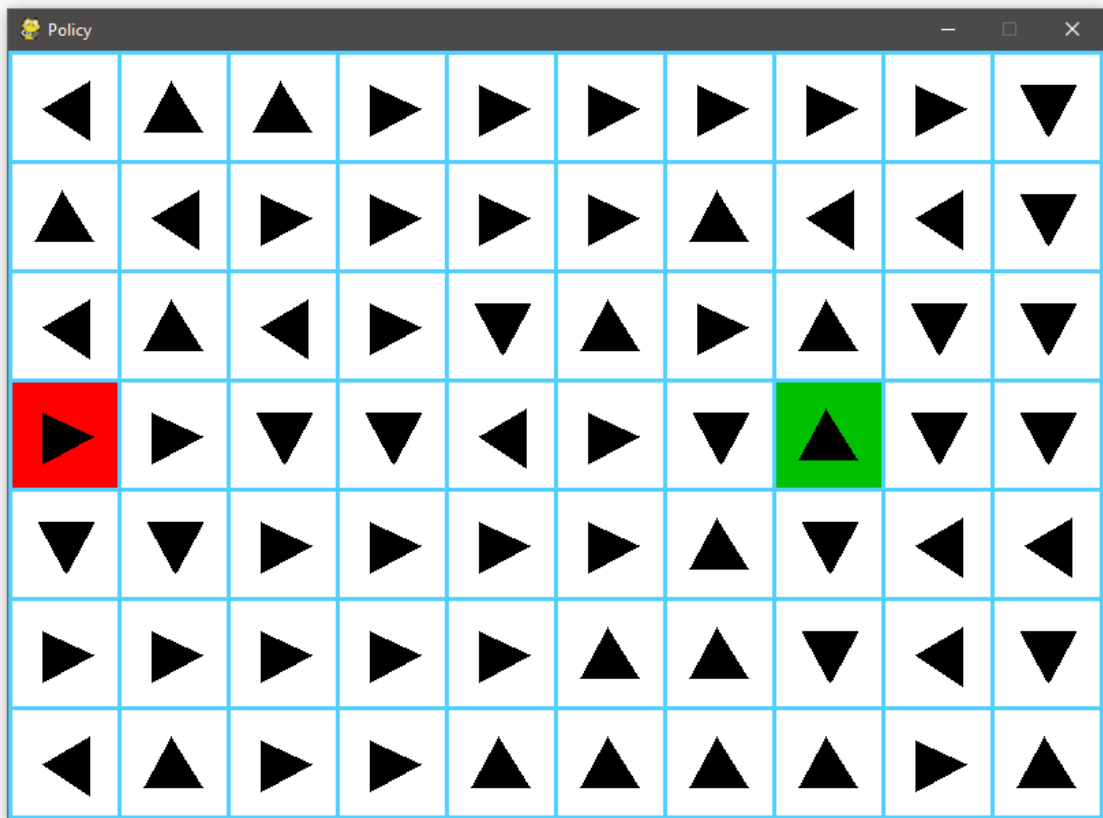
## Second: Sweeping Gamma (Discounting Rate)

- Alpha ( $\alpha$ ): 0.5
- Gamma ( $\gamma$ ): 0.2
- Epsilon ( $\epsilon$ ): 0.1
- Number of Episodes: 300

Minimal number of Episodes: 78, 55, 39, 51, 61

Average number of Episodes: 56

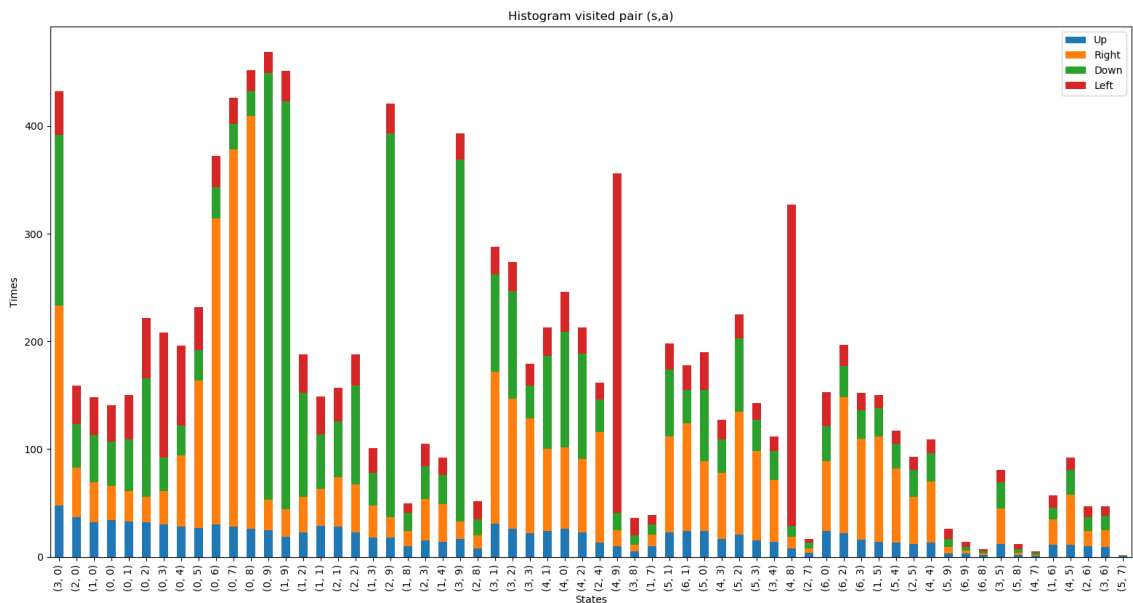
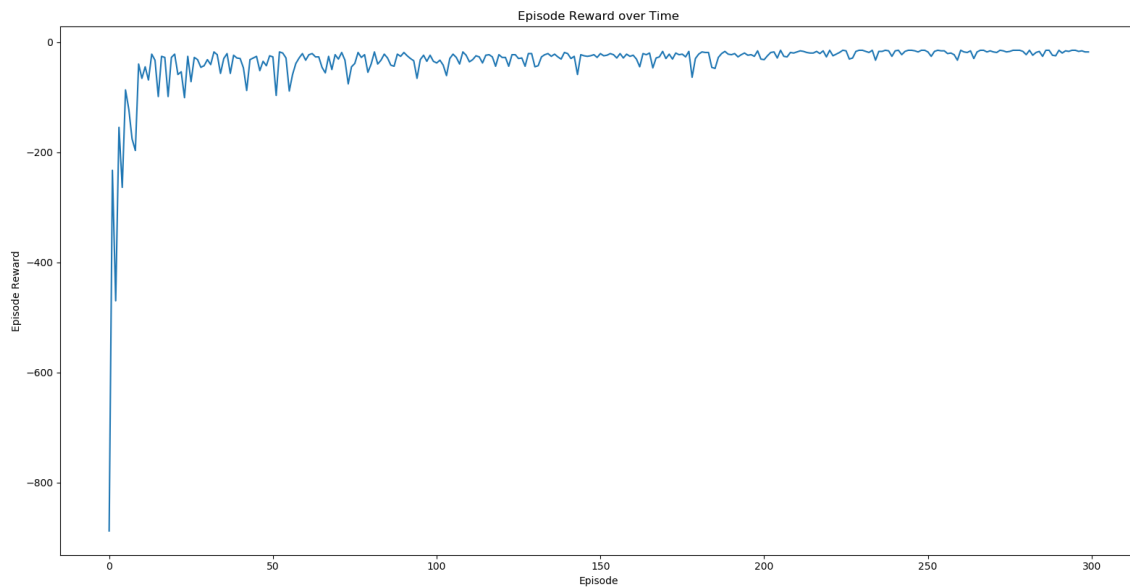




- **Alpha ( $\alpha$ ): 0.5**
- **Gamma ( $\gamma$ ): 0.5**
- **Epsilon ( $\epsilon$ ): 0.1**
- **Number of Episodes: 300**

**Minimal number of Episodes: 57, 43, 37, 45, 46**

**Average number of Episodes: 46**

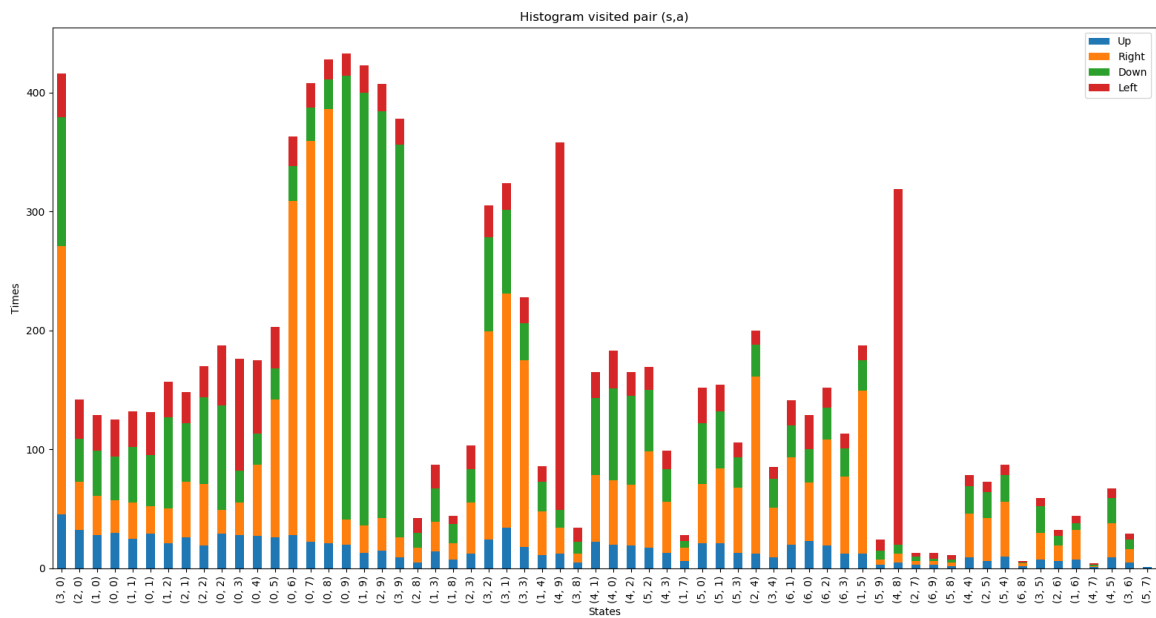
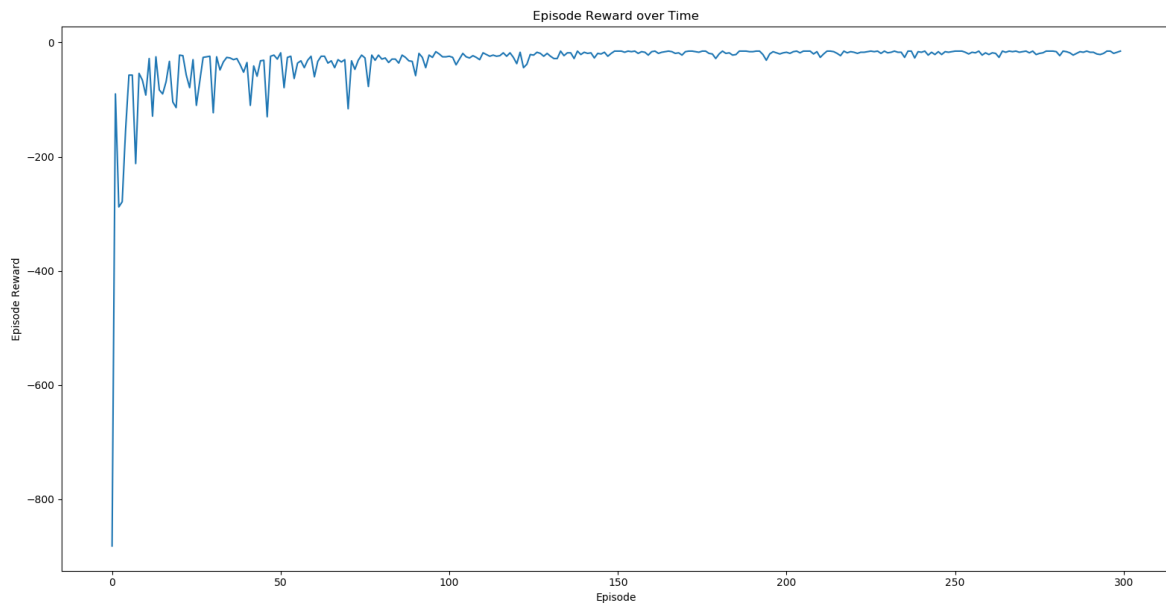




- Alpha ( $\alpha$ ): 0.5
- Gamma ( $\gamma$ ): 0.8
- Epsilon ( $\epsilon$ ): 0.1
- Number of Episodes: 300

Minimal number of Episodes: 72, 70, 102, 79, 70

Average number of Episodes: 79





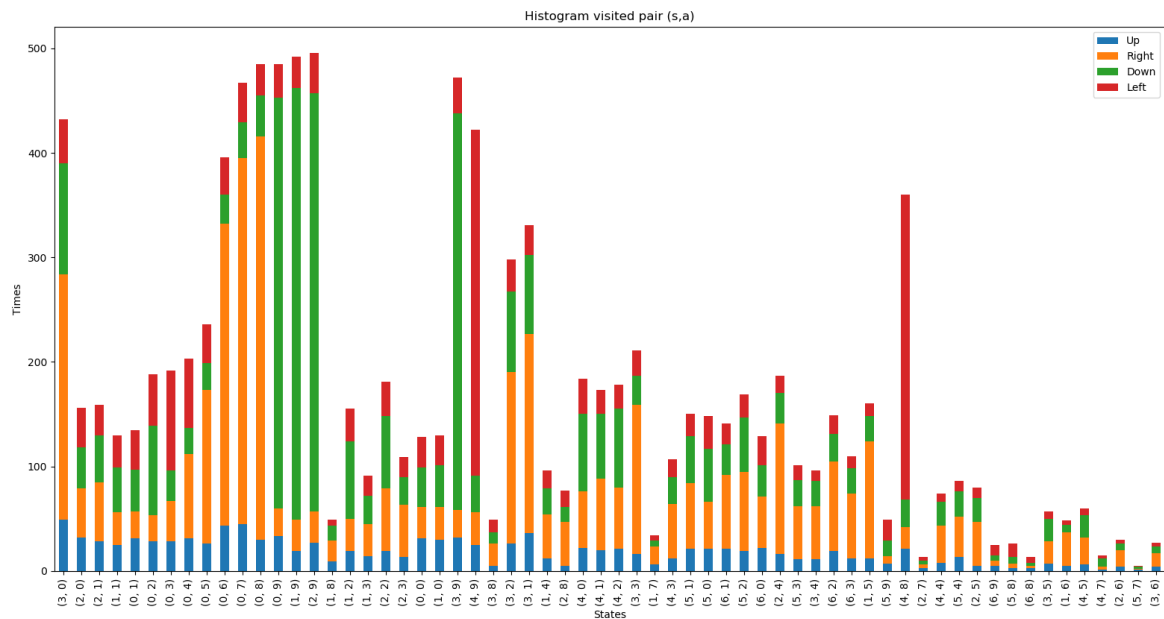
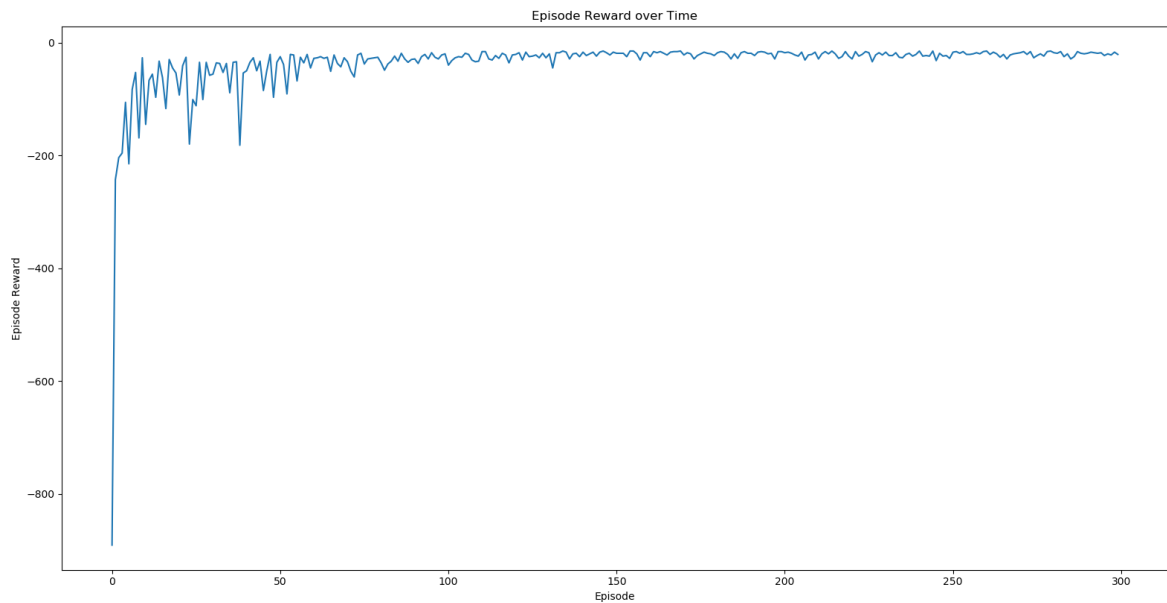


### Third: Sweeping Epsilon (Exploring Rate)

- **Alpha ( $\alpha$ ):** 0.5
- **Gamma ( $\gamma$ ):** 1
- **Epsilon ( $\epsilon$ ):** 0.2
- **Number of Episodes:** 300

**Minimal number of Episodes:** 48, 76, 43, 66, 72

**Average number of Episodes:** 61

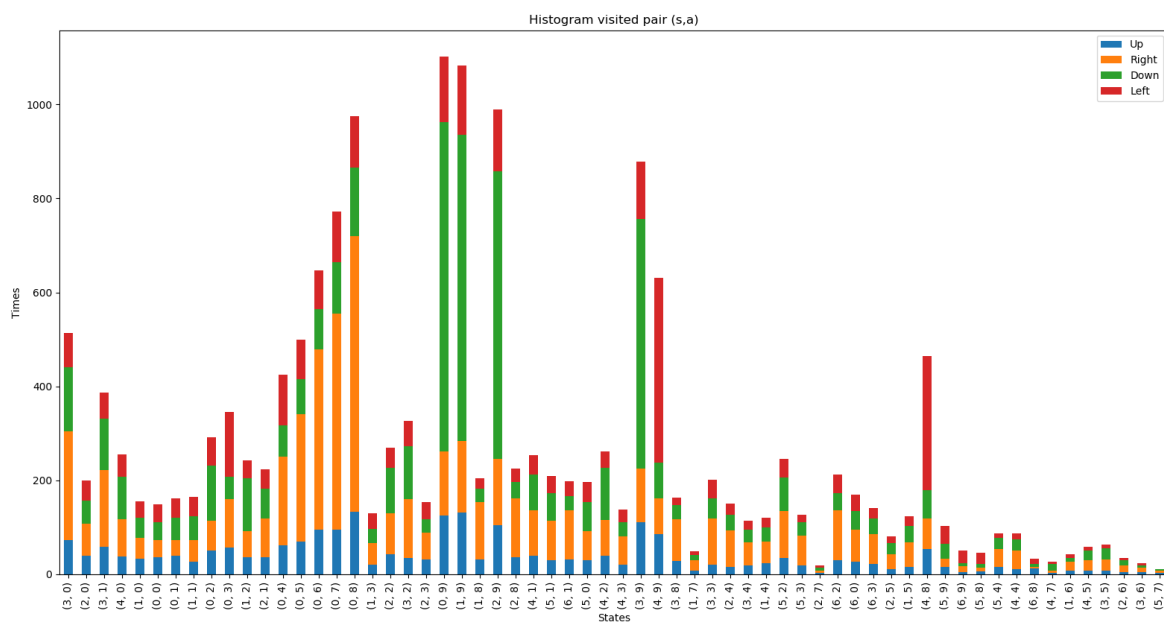
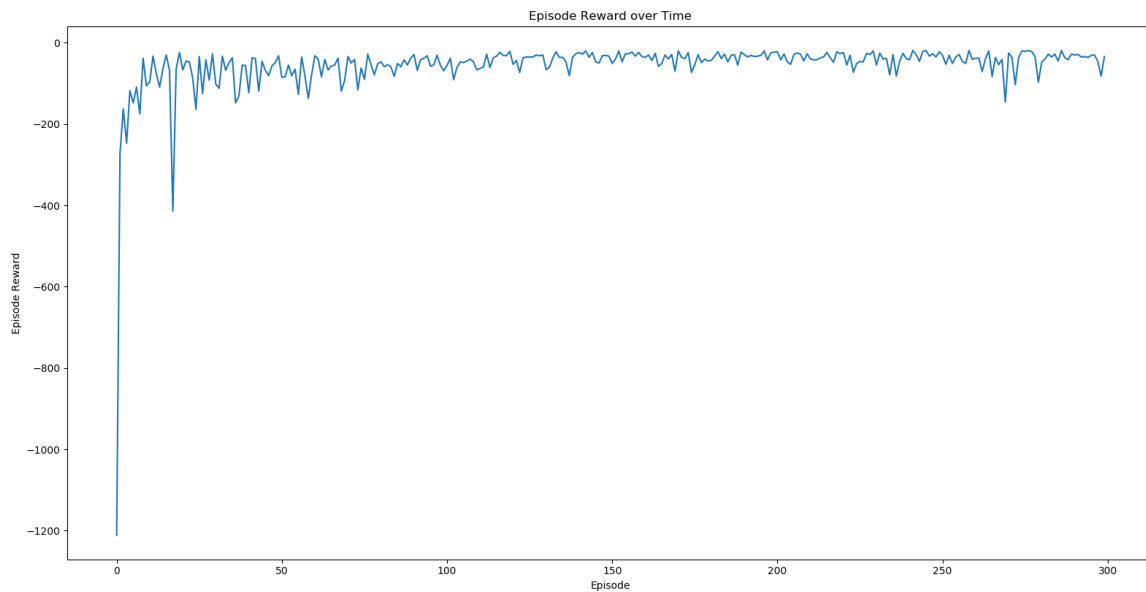




- Alpha ( $\alpha$ ): 0.5
- Gamma ( $\gamma$ ): 1
- Epsilon ( $\epsilon$ ): 0.5
- Number of Episodes: 300

Minimal number of Episodes: 37, 31, 30, 32, 36

Average number of Episodes: 33

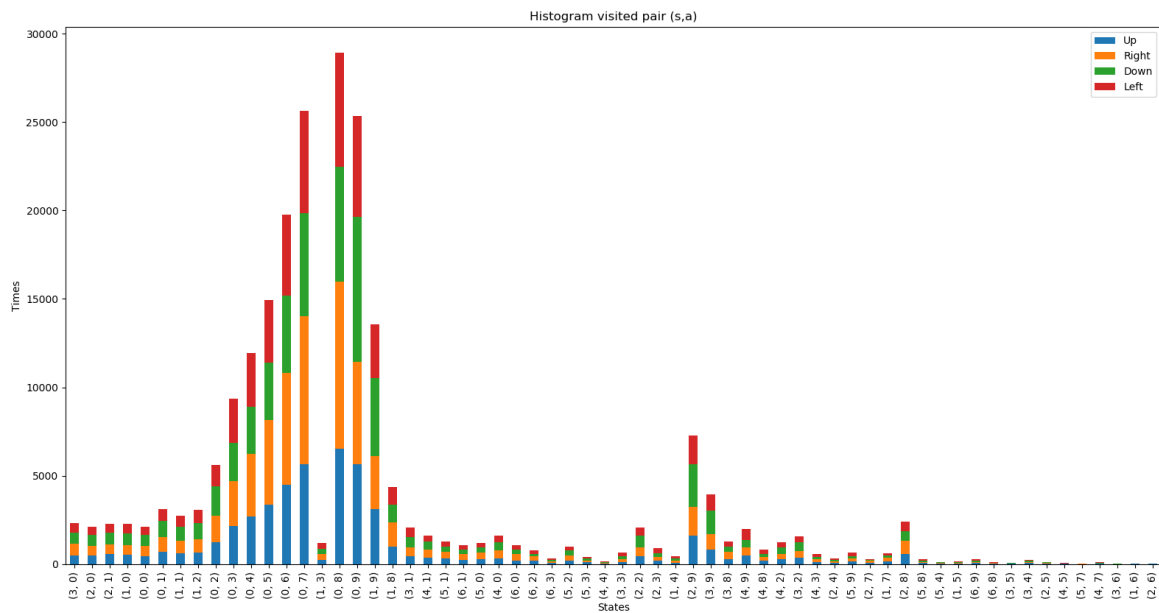
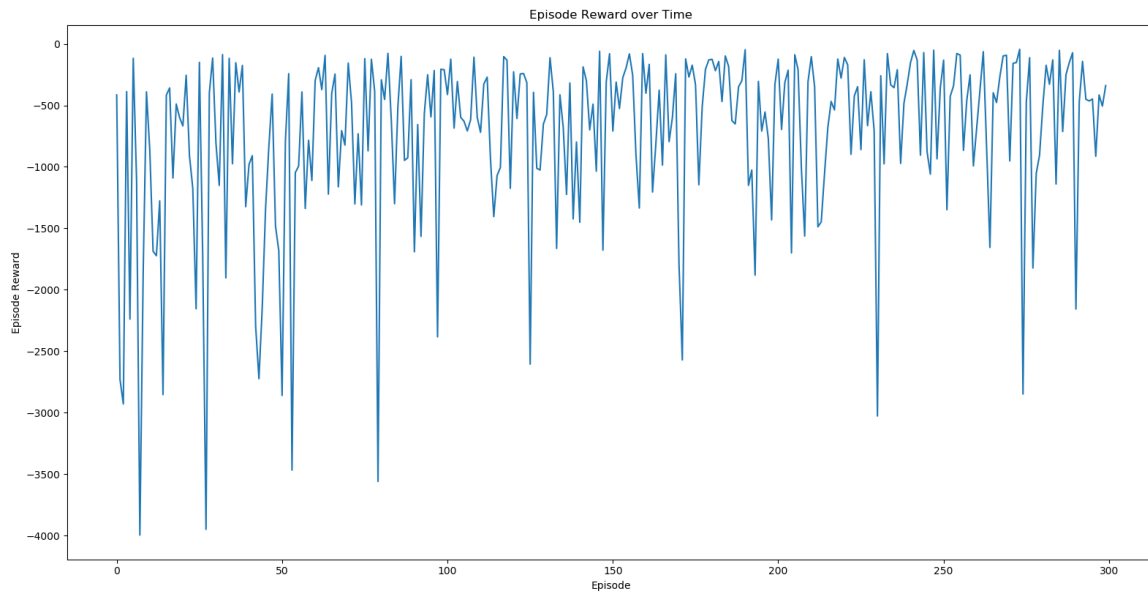




- Alpha ( $\alpha$ ): 0.5
- Gamma ( $\gamma$ ): 1
- Epsilon ( $\epsilon$ ): 0.9
- Number of Episodes: 300

Minimal number of Episodes: 129, 416, 276, 432, 100

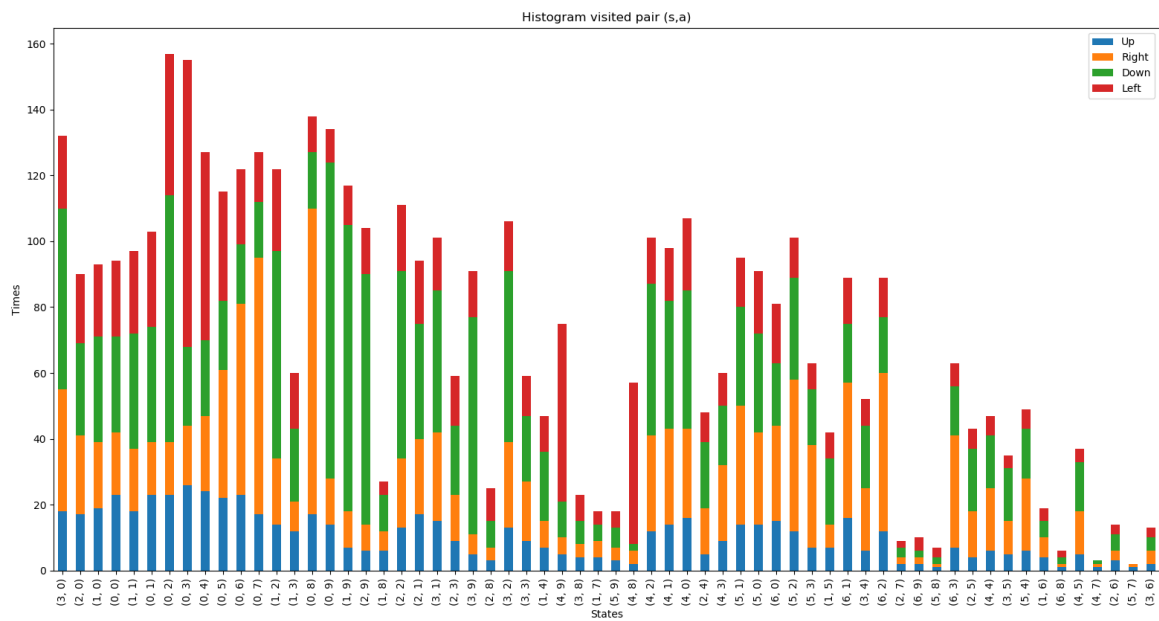
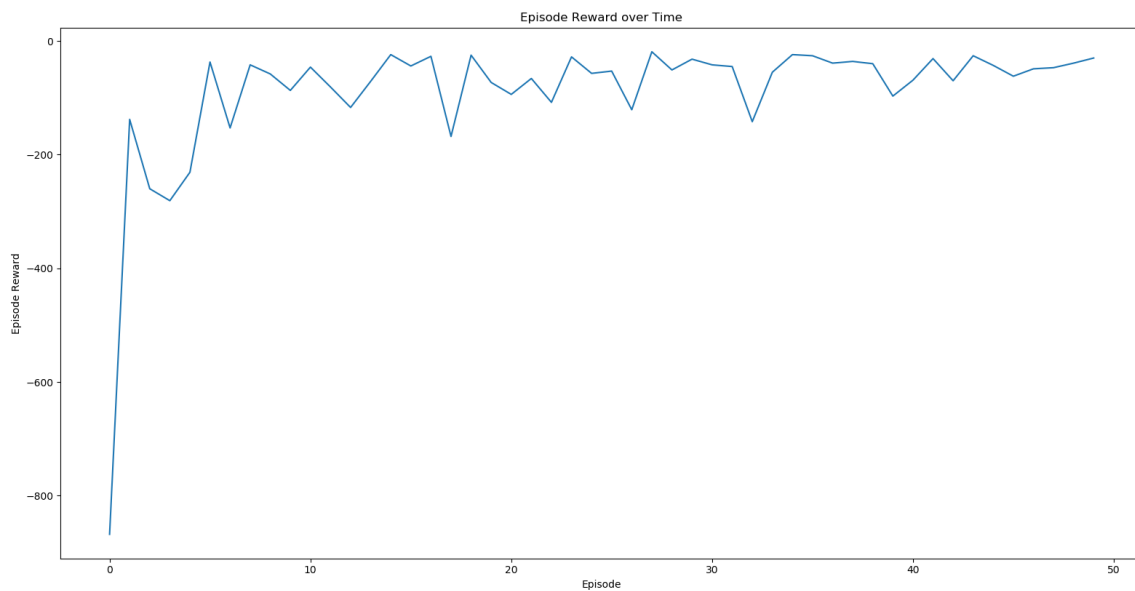
Average number of Episodes: 271

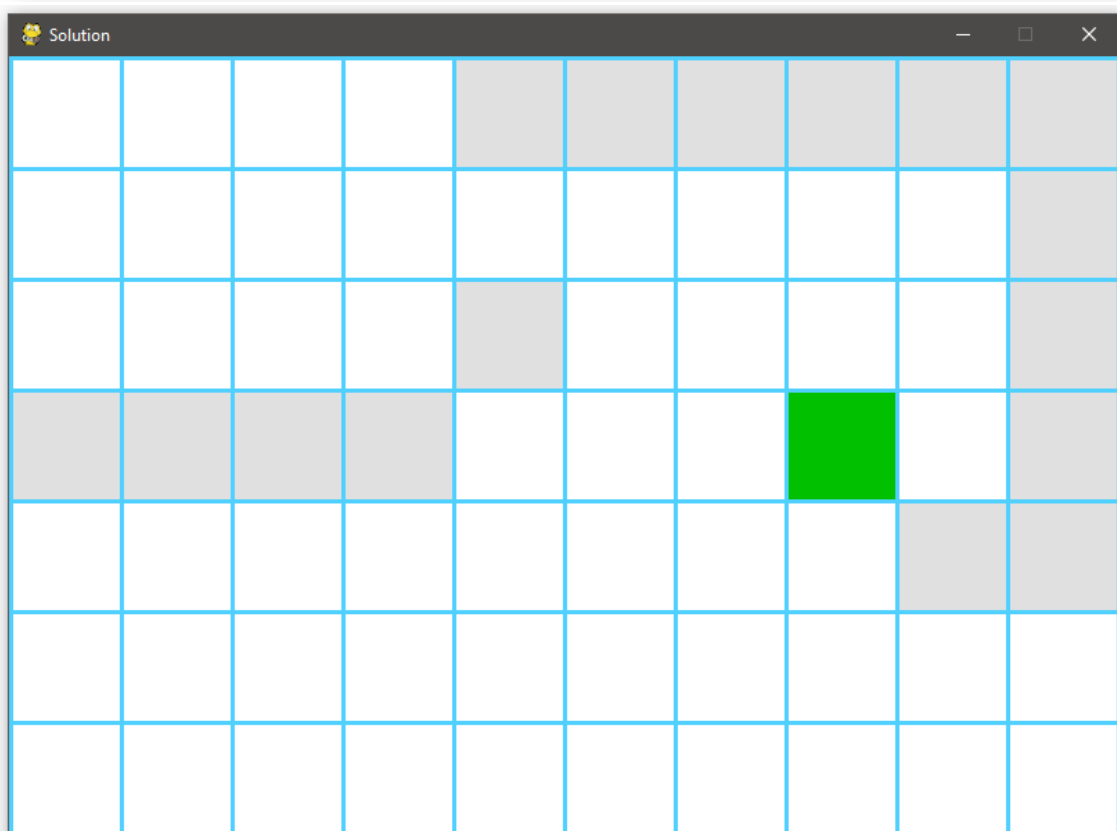
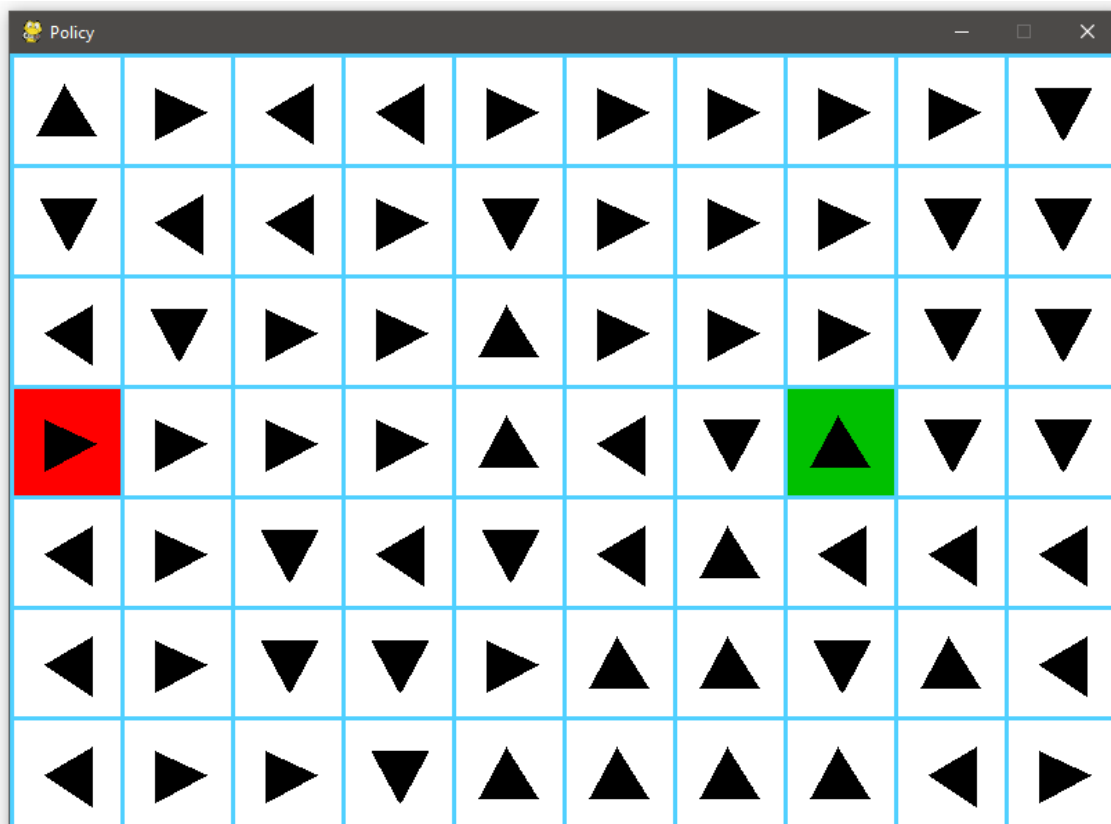




#### Fourth: Sweeping Episodes

- **Alpha ( $\alpha$ ): 0.5**
- **Gamma ( $\gamma$ ): 1**
- **Epsilon ( $\epsilon$ ): 0.1**
- **Number of Episodes: 50**

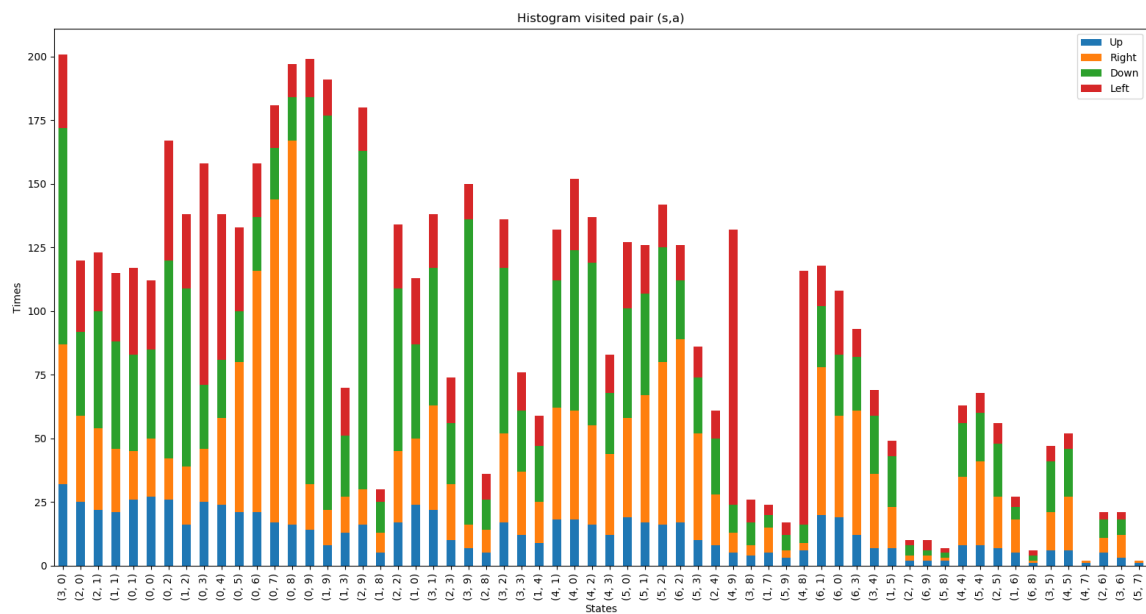
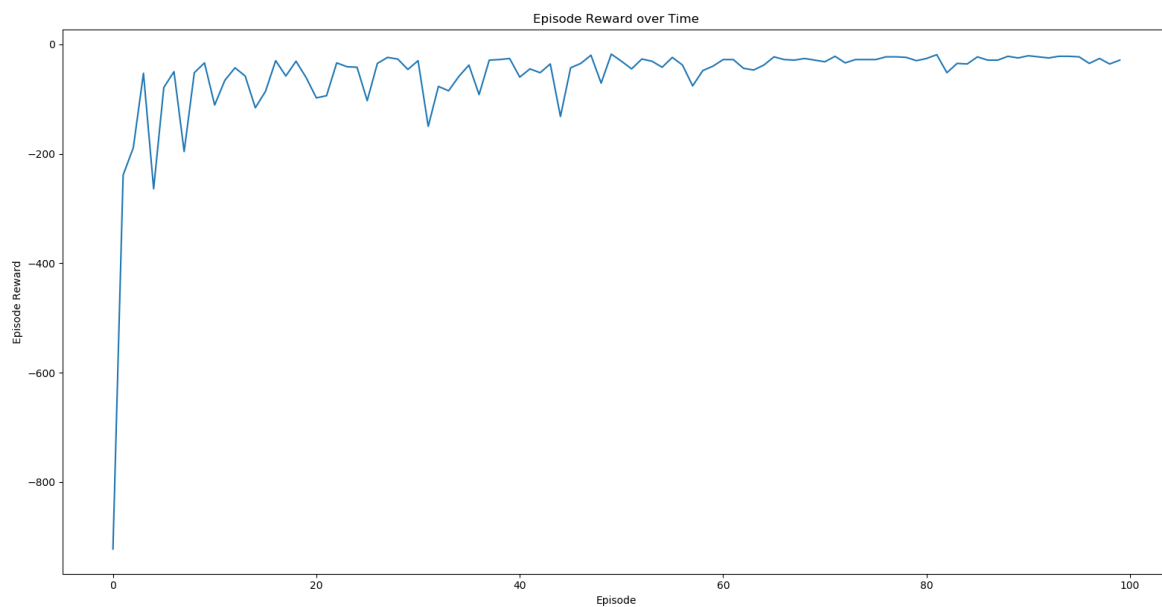






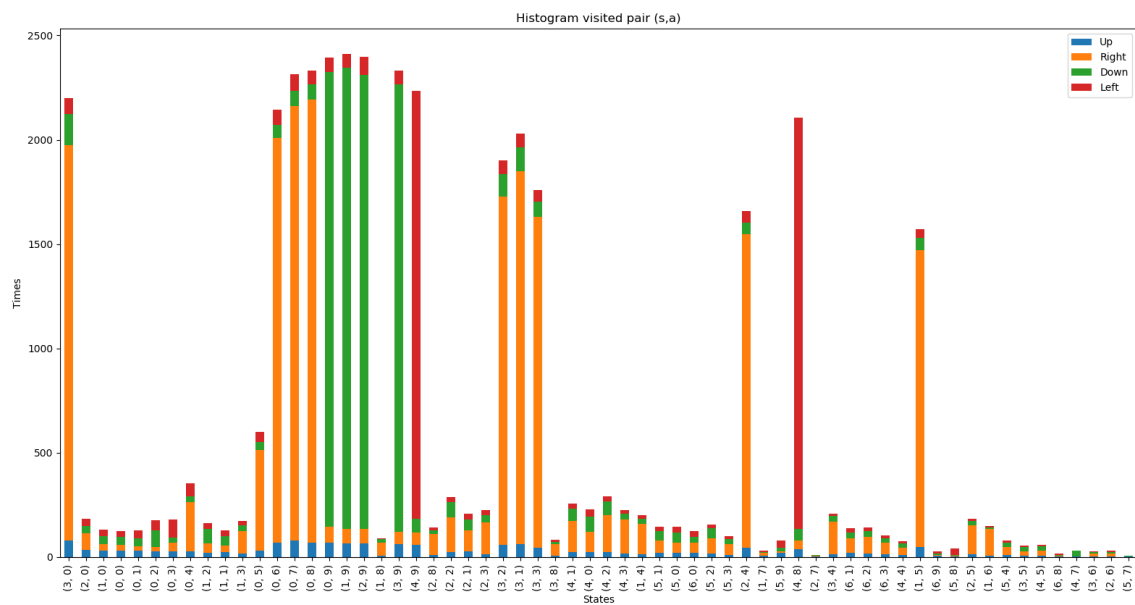
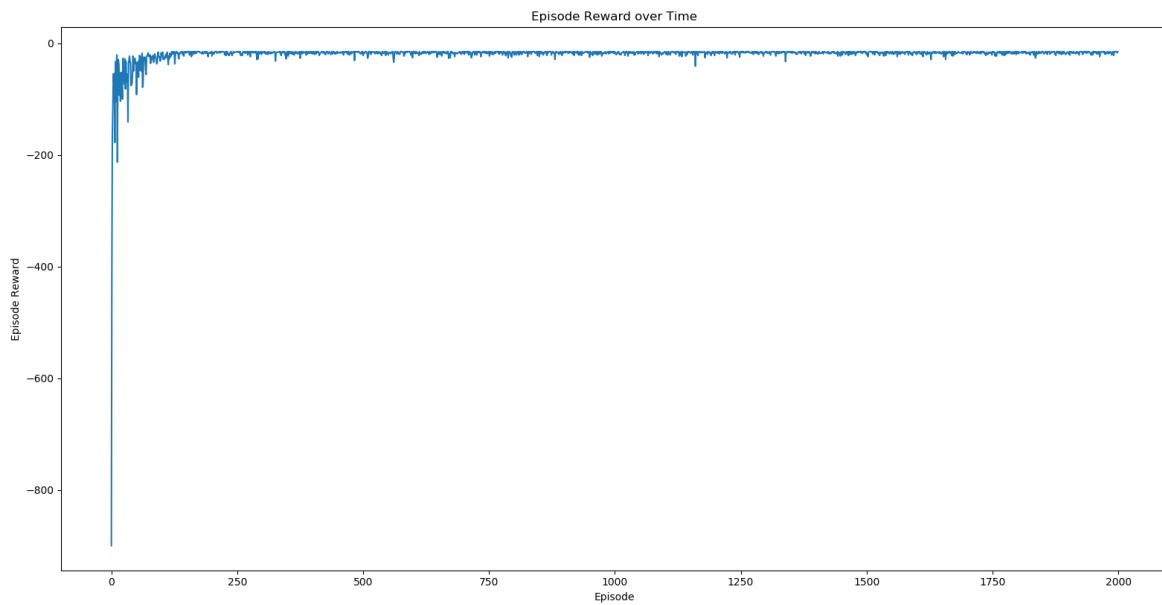
**Alpha ( $\alpha$ ): 0.5**

- **Gamma ( $\gamma$ ): 1**
- **Epsilon ( $\epsilon$ ): 0.1**
- **Number of Episodes: 100**





- **Alpha ( $\alpha$ ):** 0.5
- **Gamma ( $\gamma$ ):** 1
- **Epsilon ( $\epsilon$ ):** 0.1
- **Number of Episodes:** 2000





Just looking at the average episodes needed to get a path, we can say that the best configuration is:

- **Alpha ( $\alpha$ ):** 0.5
- **Gamma ( $\gamma$ ):** 1
- **Epsilon ( $\epsilon$ ):** 0.5
- **Number of Episodes:** 300

And looking through all the graphics for every change made in the variables, we can say that:

- A learning rate higher than 0.6 doesn't make sense for the current environment.
- A higher gamma doesn't mean better results, and values around 0.5 could work really good.
- The Epsilon for the greedy algorithm implemented don't work with values higher than 0.6
- Is not needed more than 500 episodes to get the policy with a good configuration.