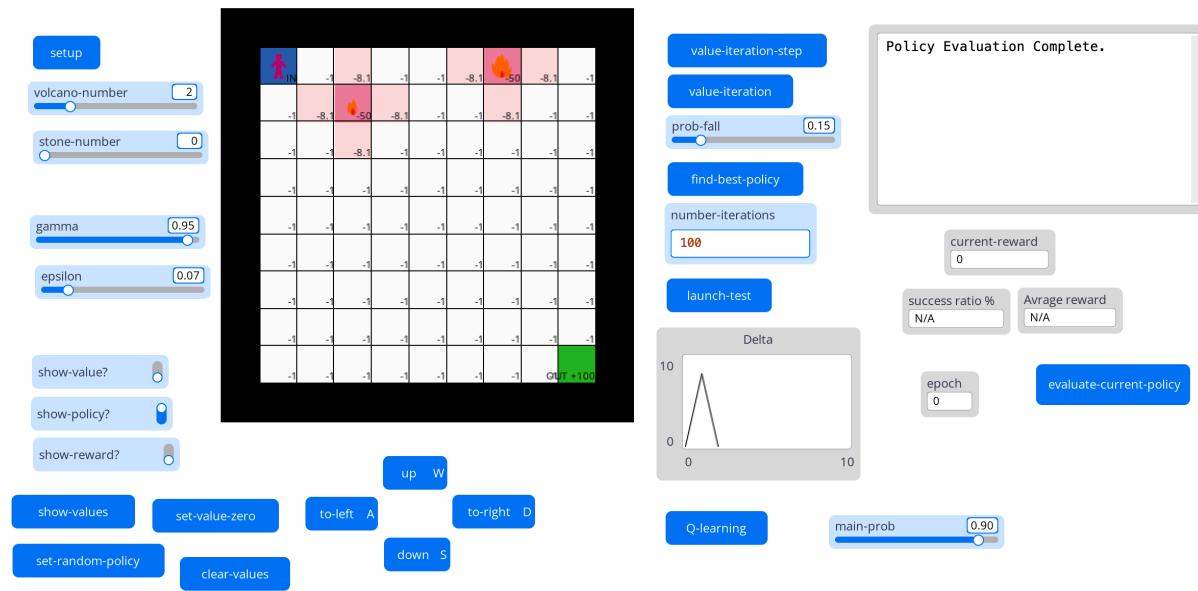


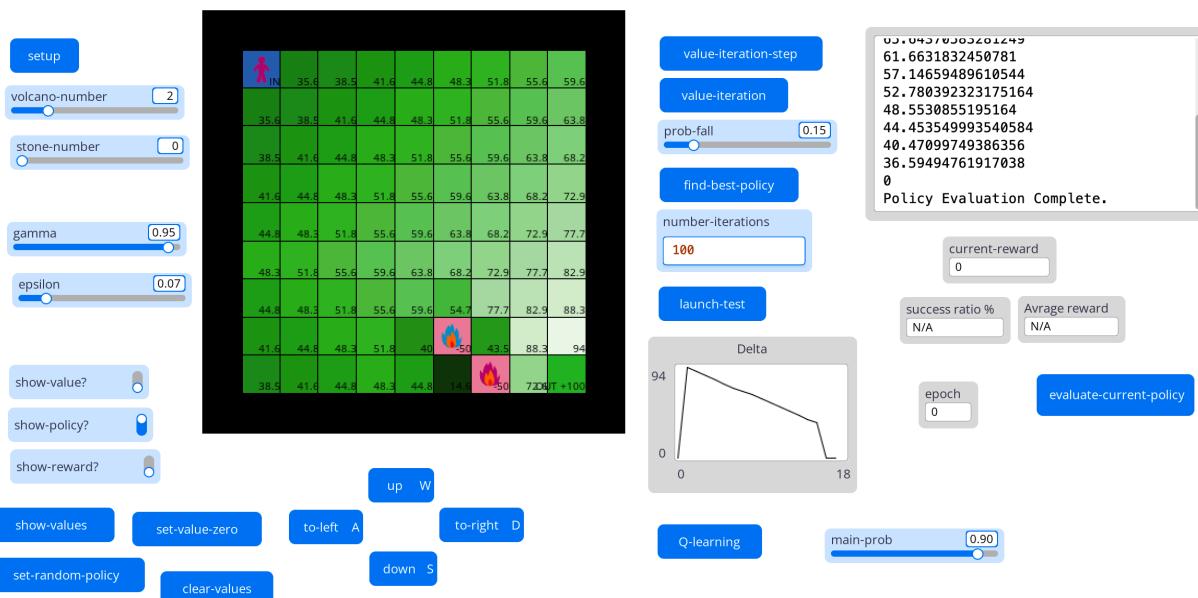
HW5 Reinforcement learning

Task 1. Policy evaluation

Random policy



Solved policy

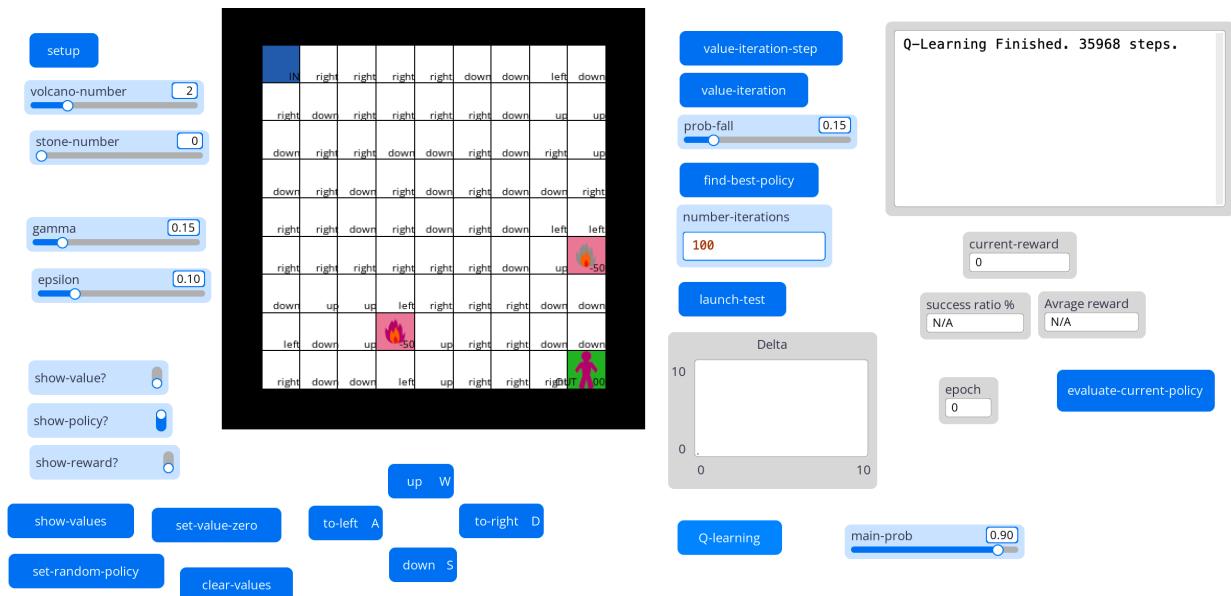


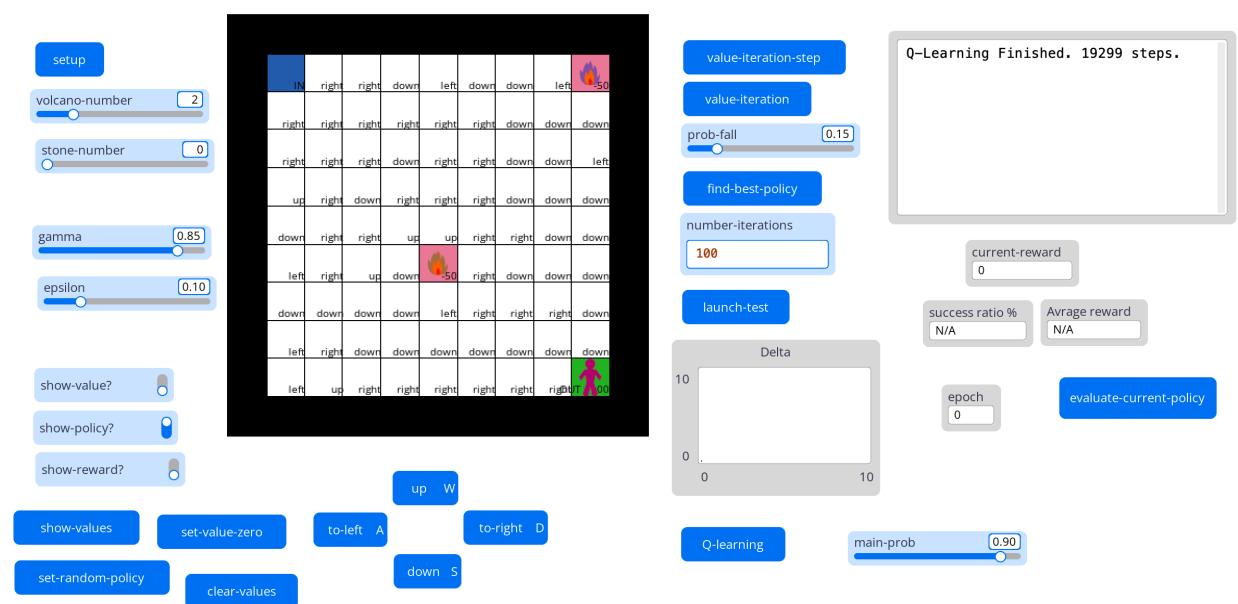
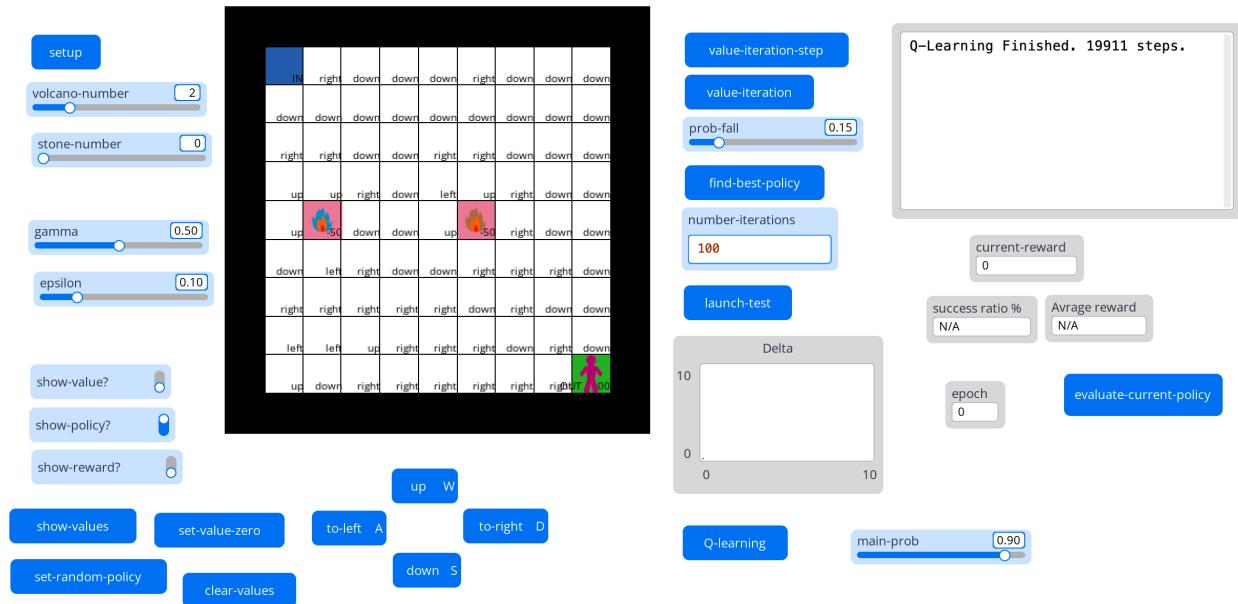
I implemented a function `evaluate-current-policy` that calculates Bellman expectation function for each cell.

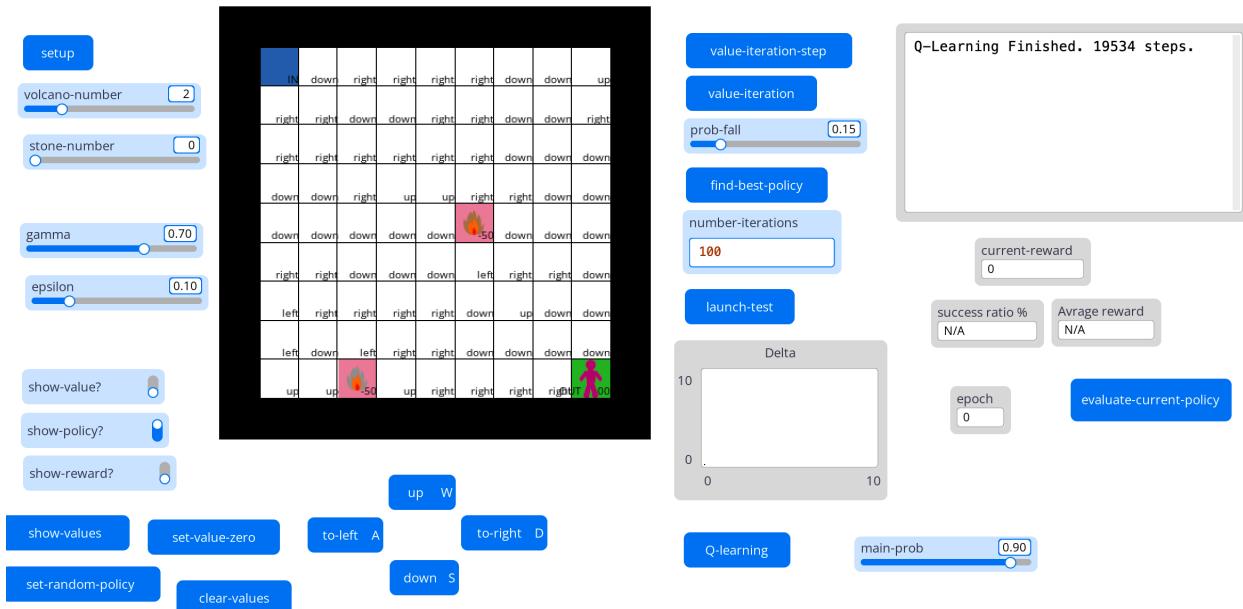
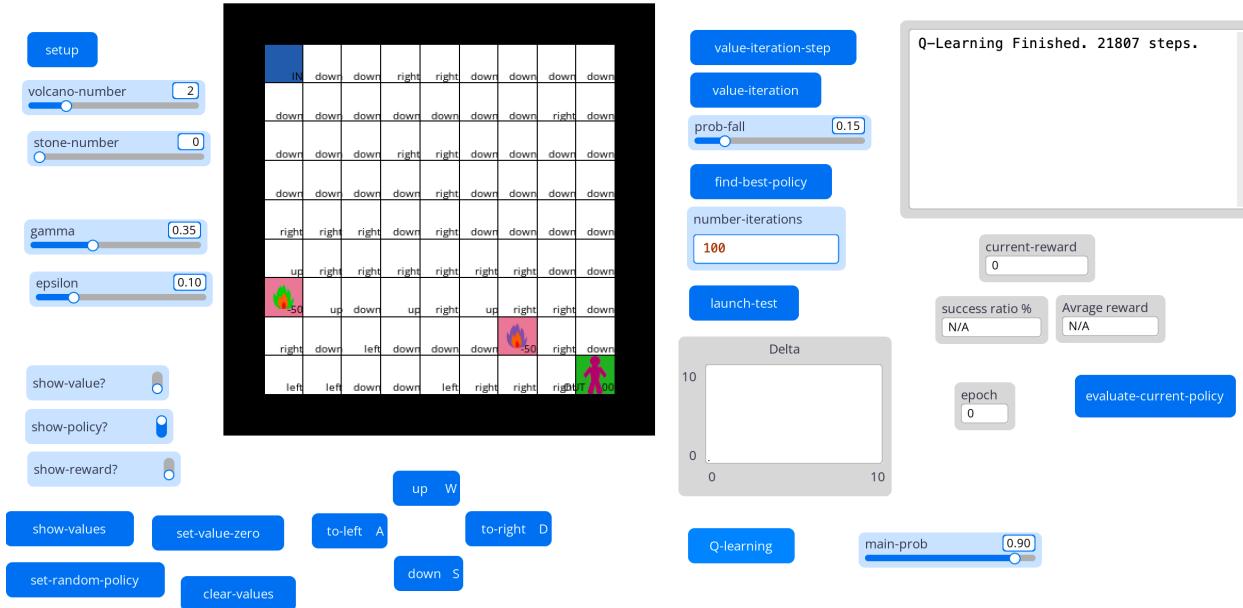
I noticed that, for random policy, cells near volcanos get lower values, which is expected.

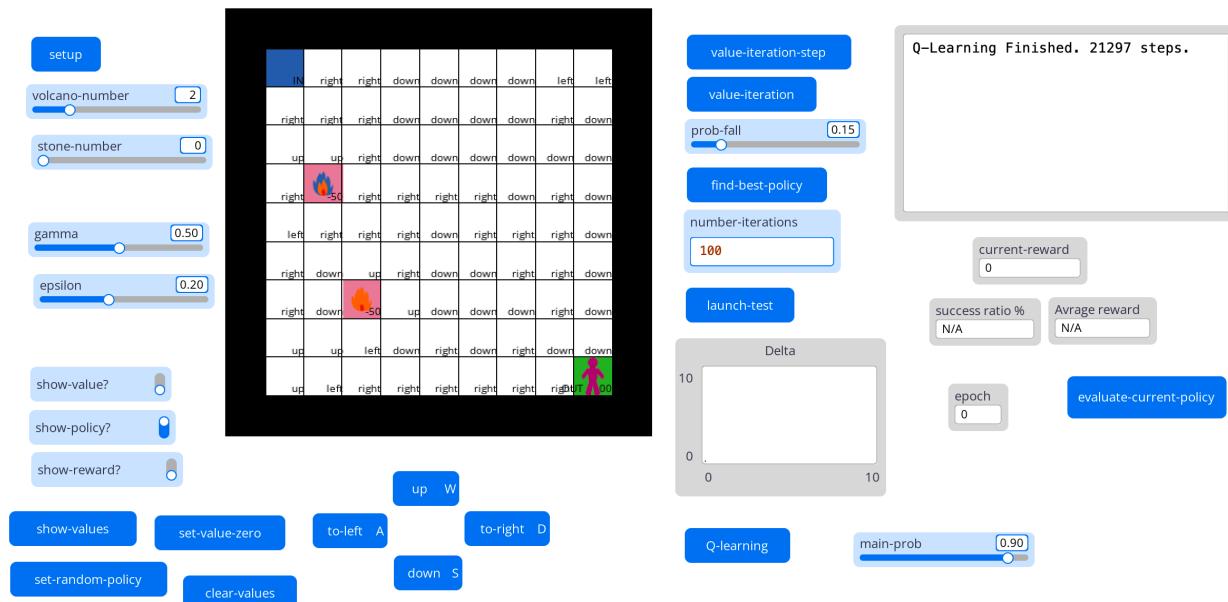
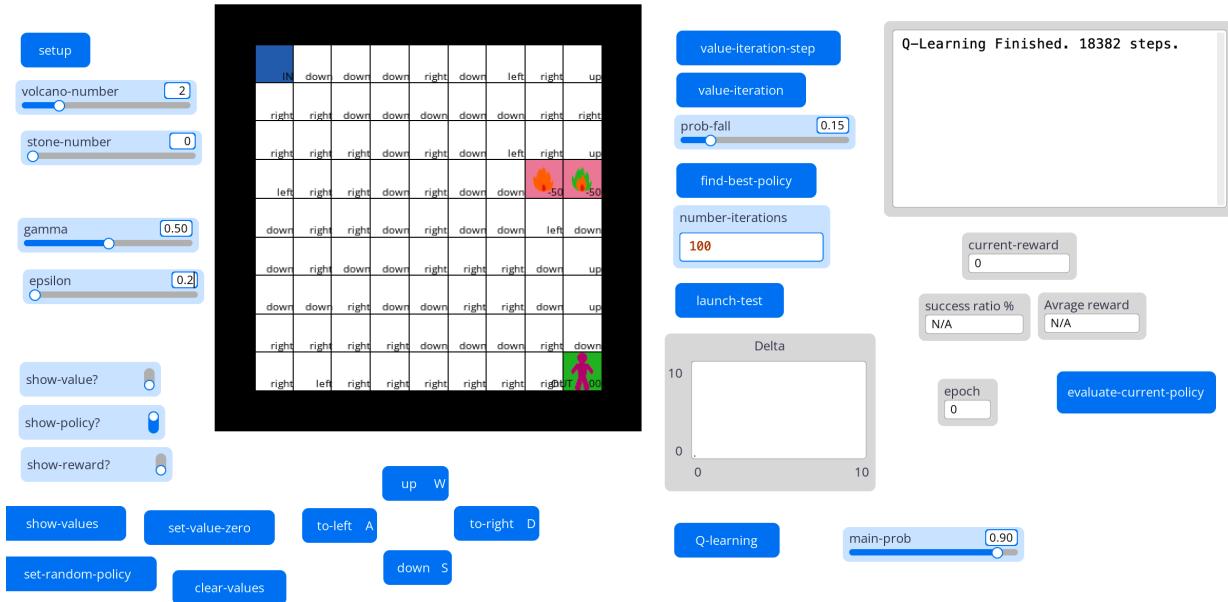
The same holds for solved policy. You can see that for solved policy almost all cells give high reward, it is explained by the fact, that there is only 2 volcanos and no obstacles, thus almost all routes lead to END.

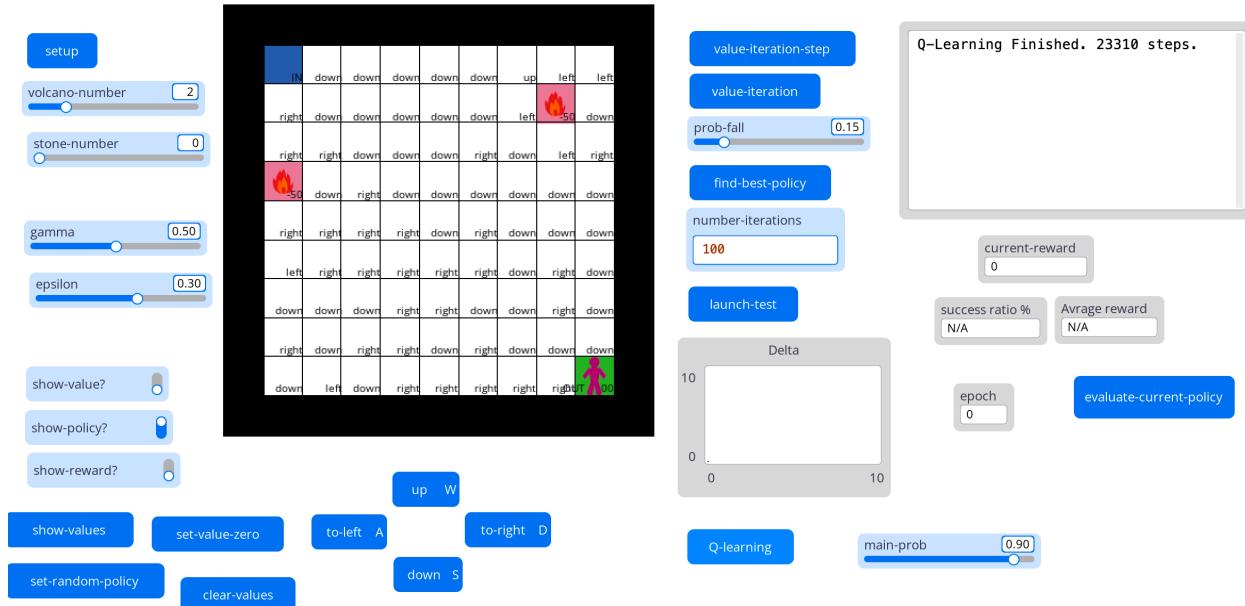
Q-learning

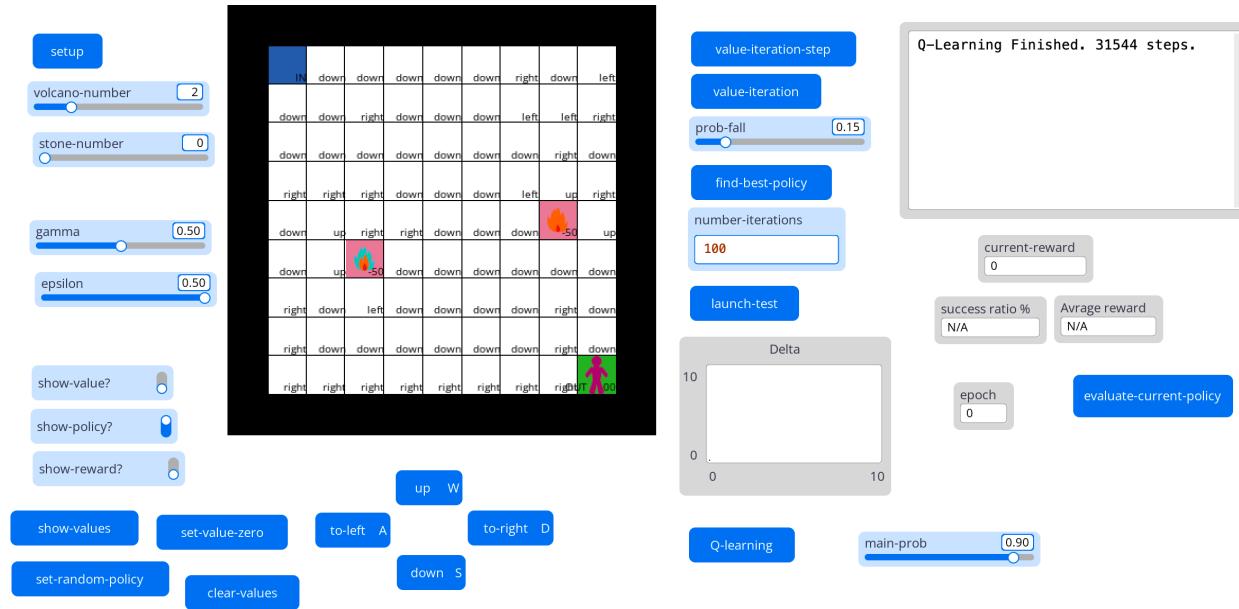












It tried different different learning rates γ and different e-e ratio epsilon. Turns out that the quickest convergence is achieved at $\gamma = 0.5, \varepsilon = 0.3$.

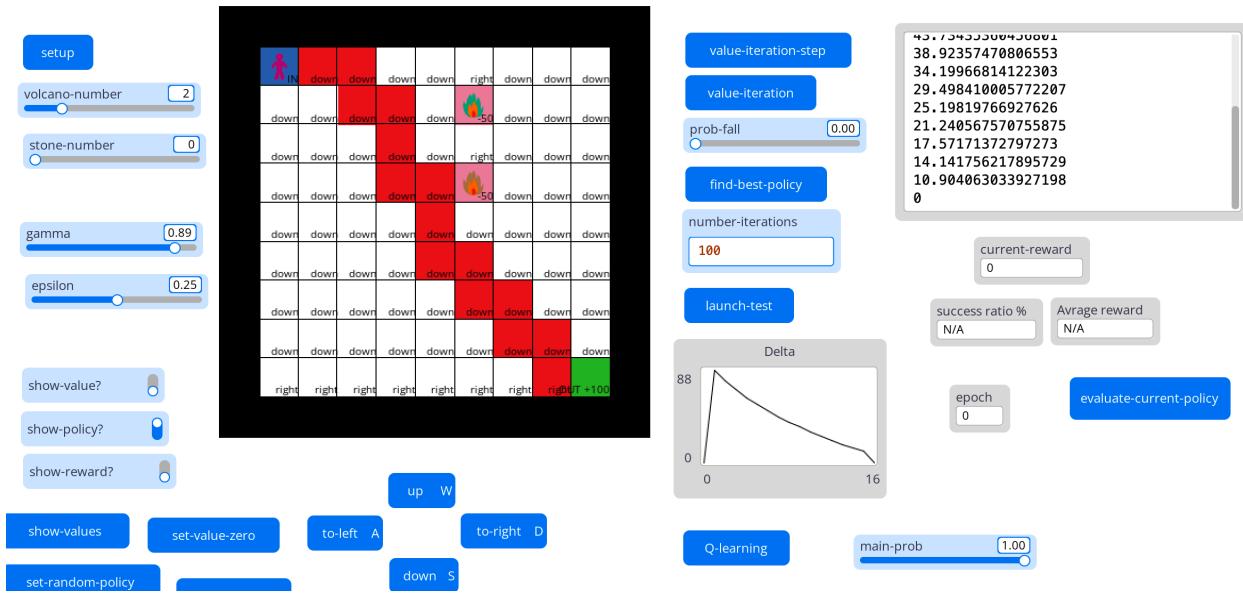


Task after Q-learning

Take MazeMDPDeterministic and add current change in value function plot and find-best-policy button.

How is it working?

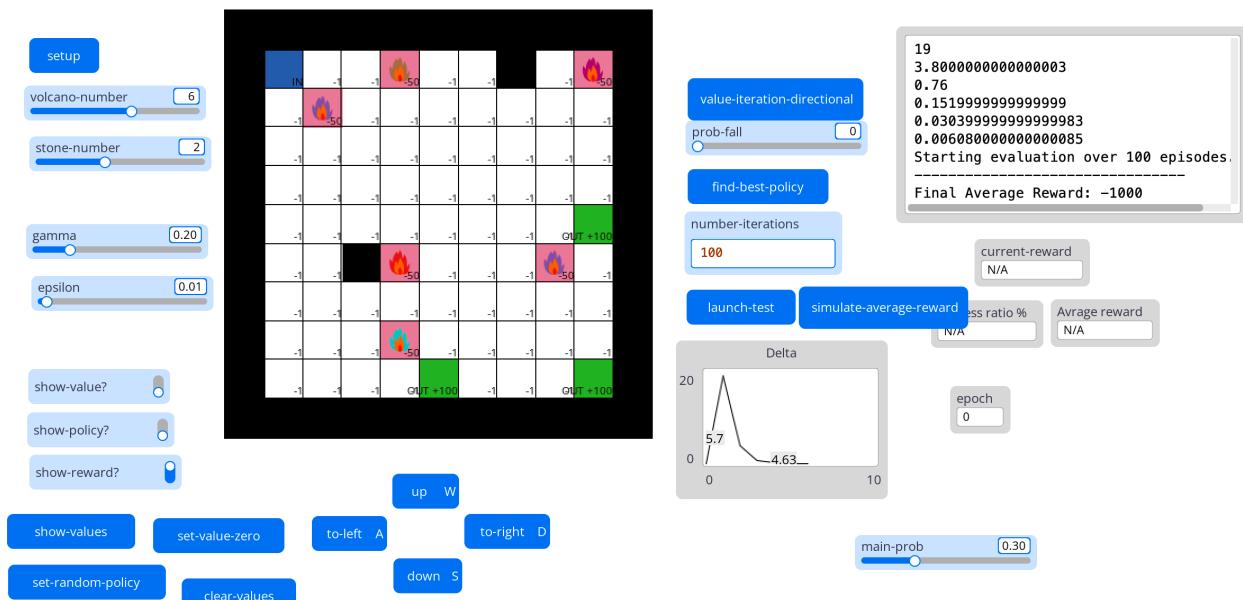
Works well. Converges

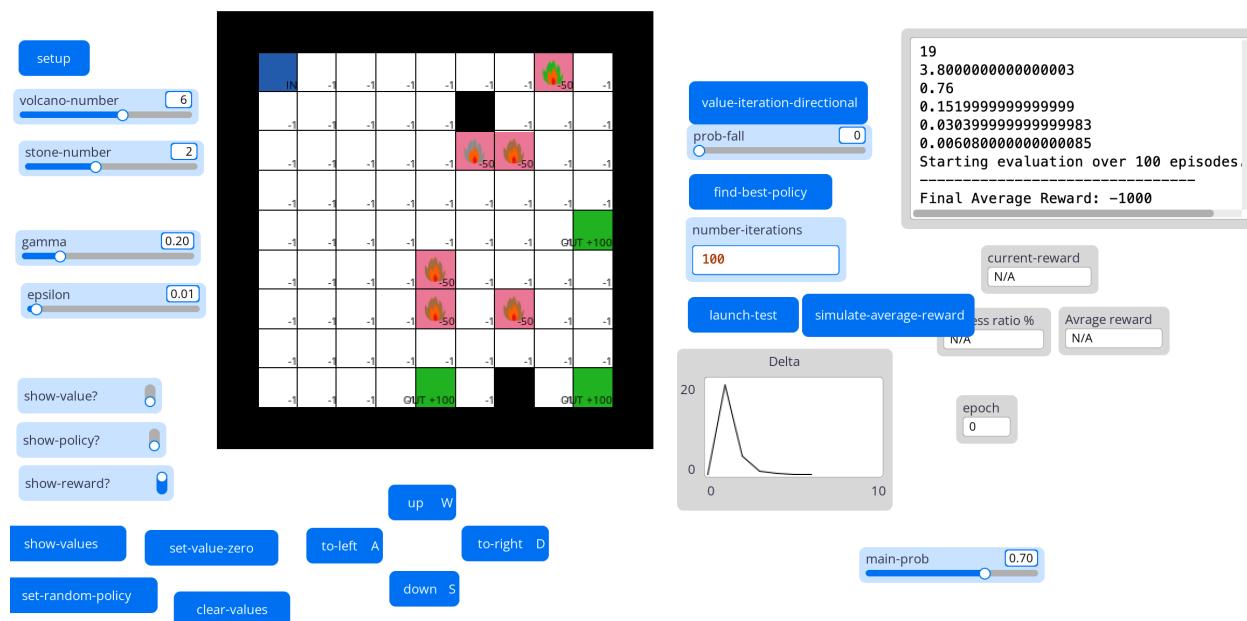
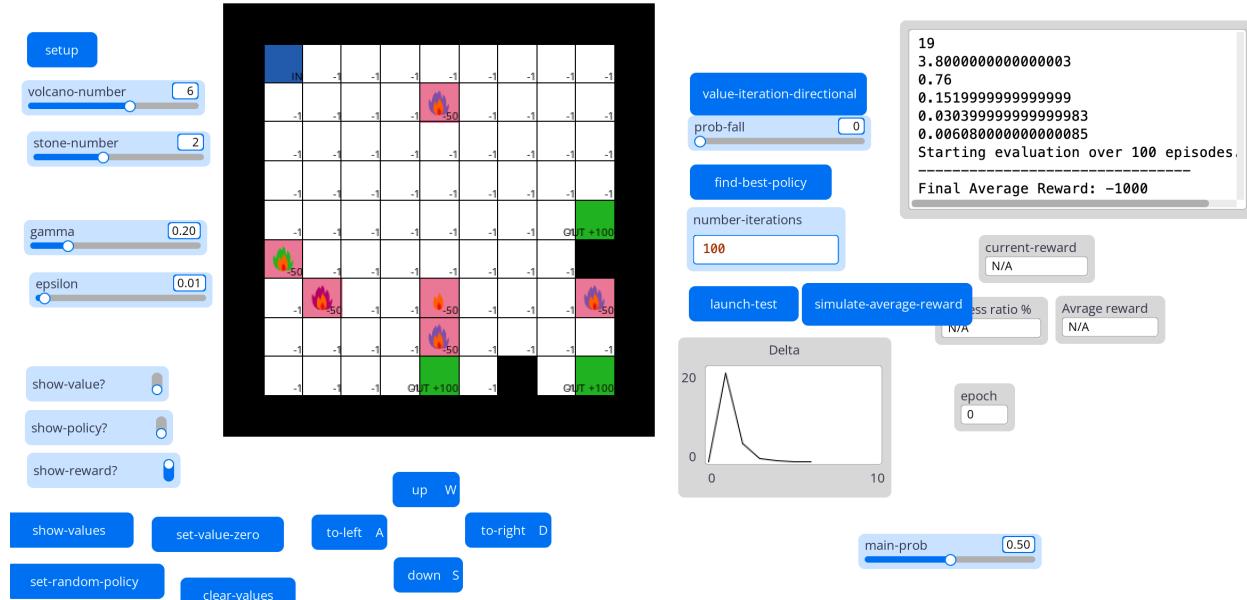


Task probabilistic case (new action model)

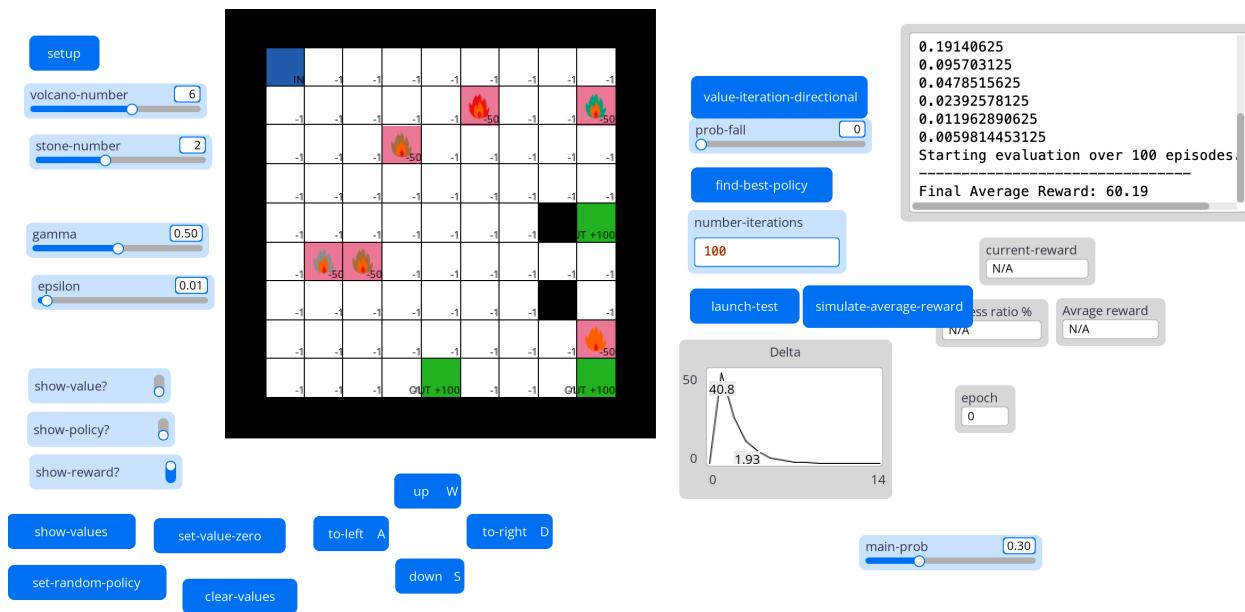
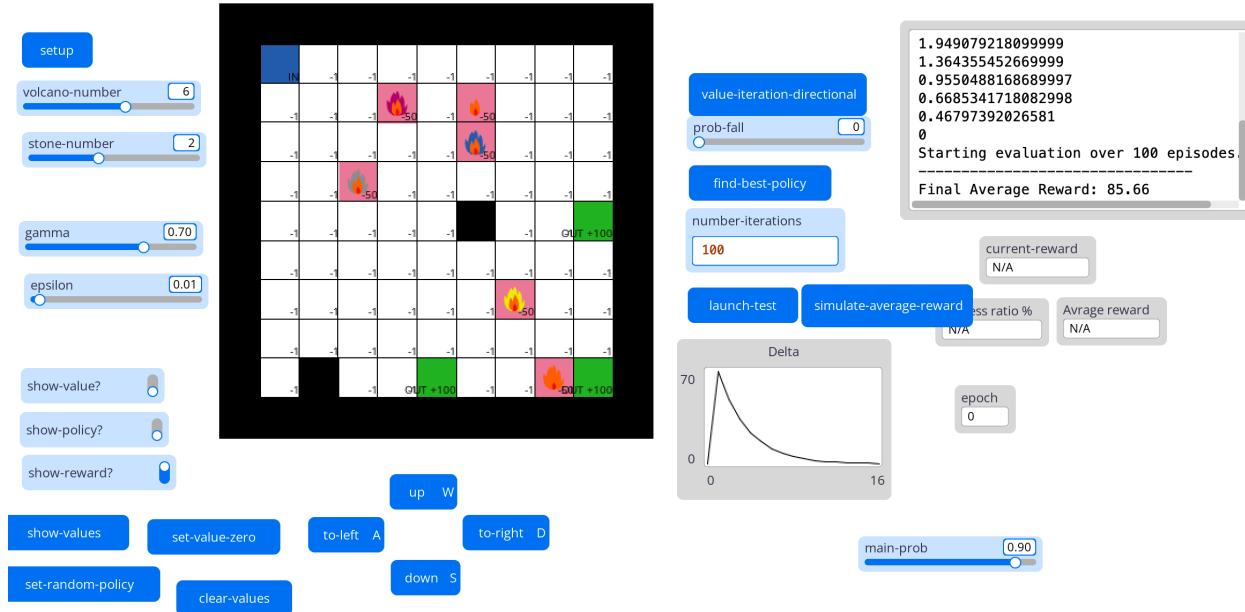
The state are now coordinate + heading (one out of four possible). Actions are now rotation to the left/right or step, each having reward (or rather punishment) -1 . To implement this model I also need to store list of 4 values for each possible direction.

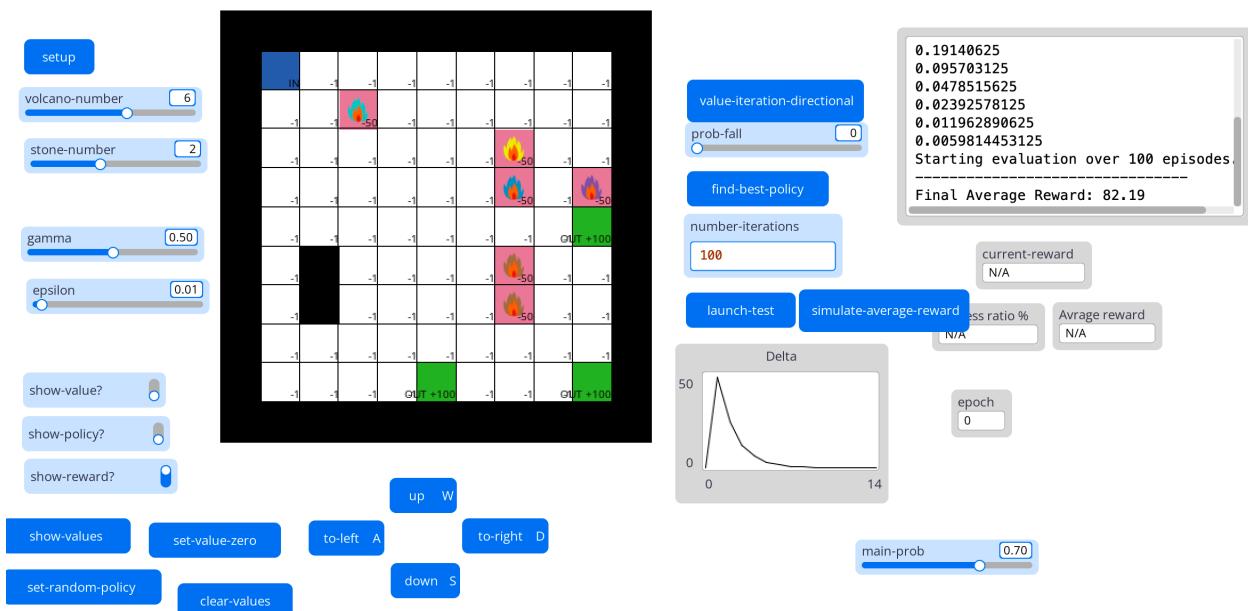
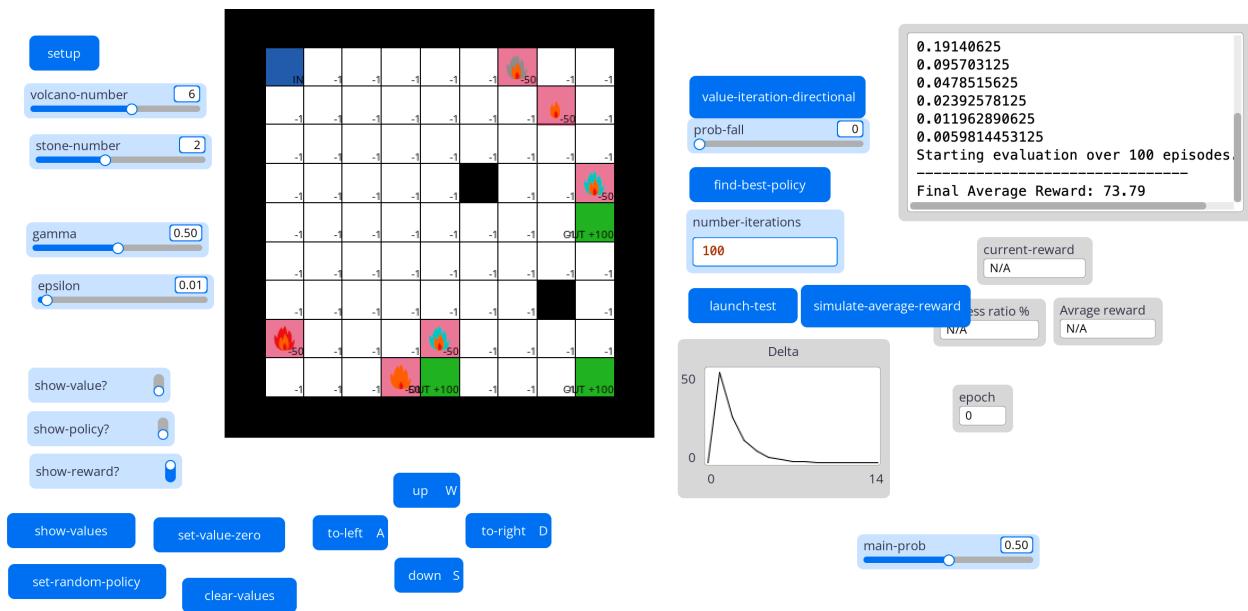
I created a function that runs solved policy for n -times and calculated average total reward

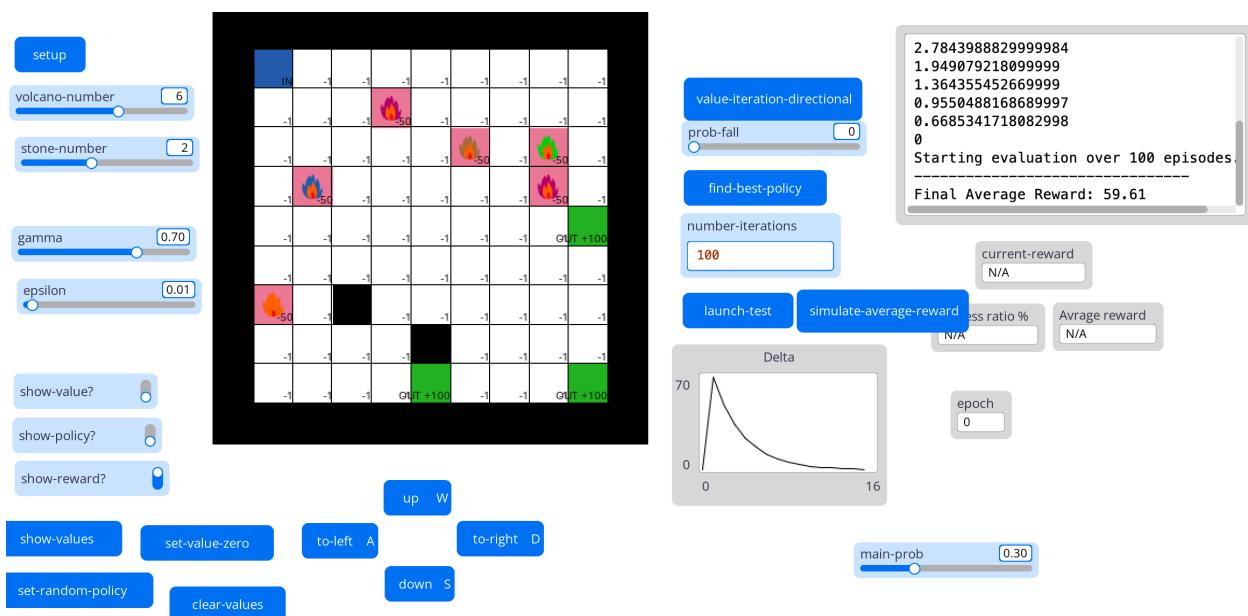
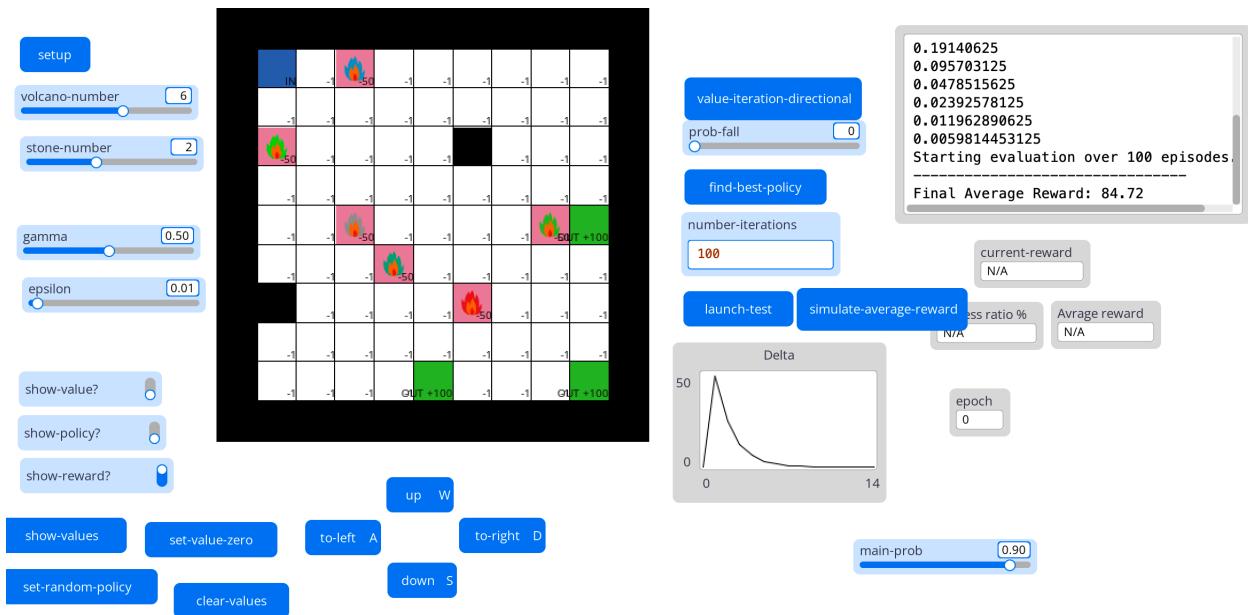


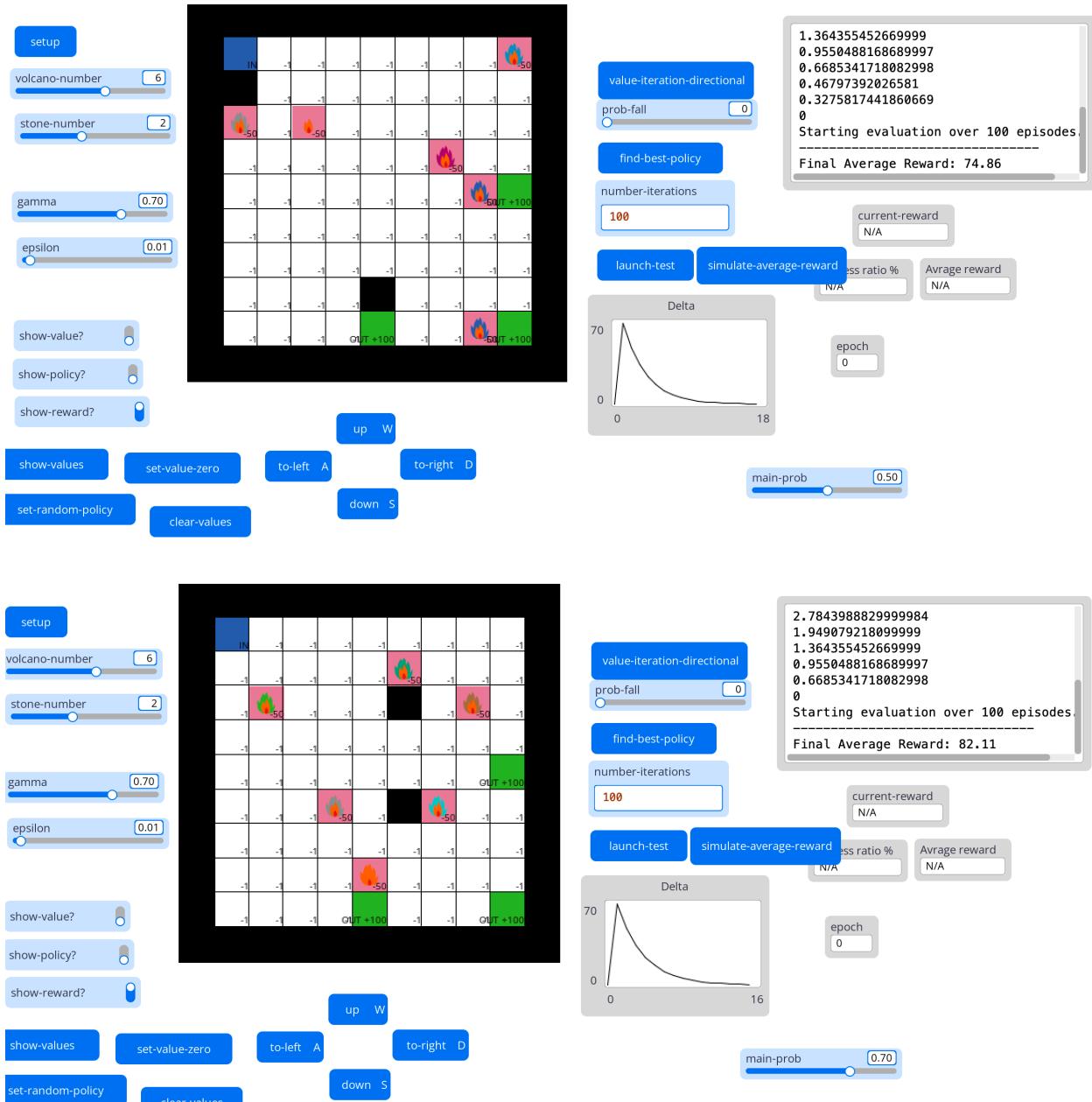


We can notice, that for low learning rate γ policy cannot be solved



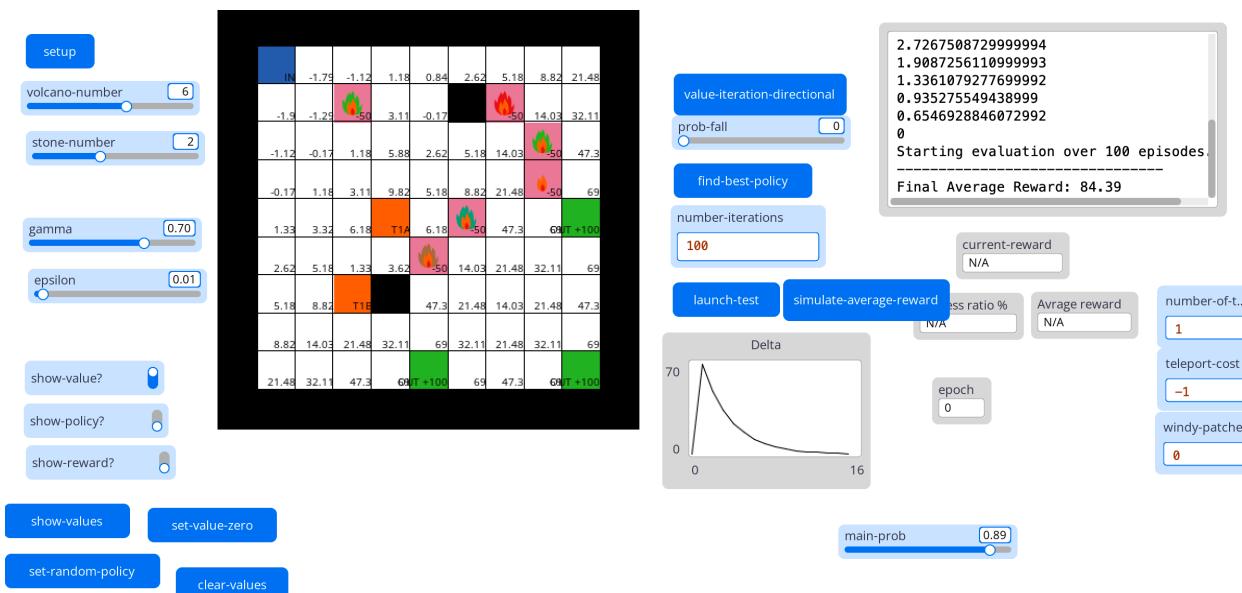
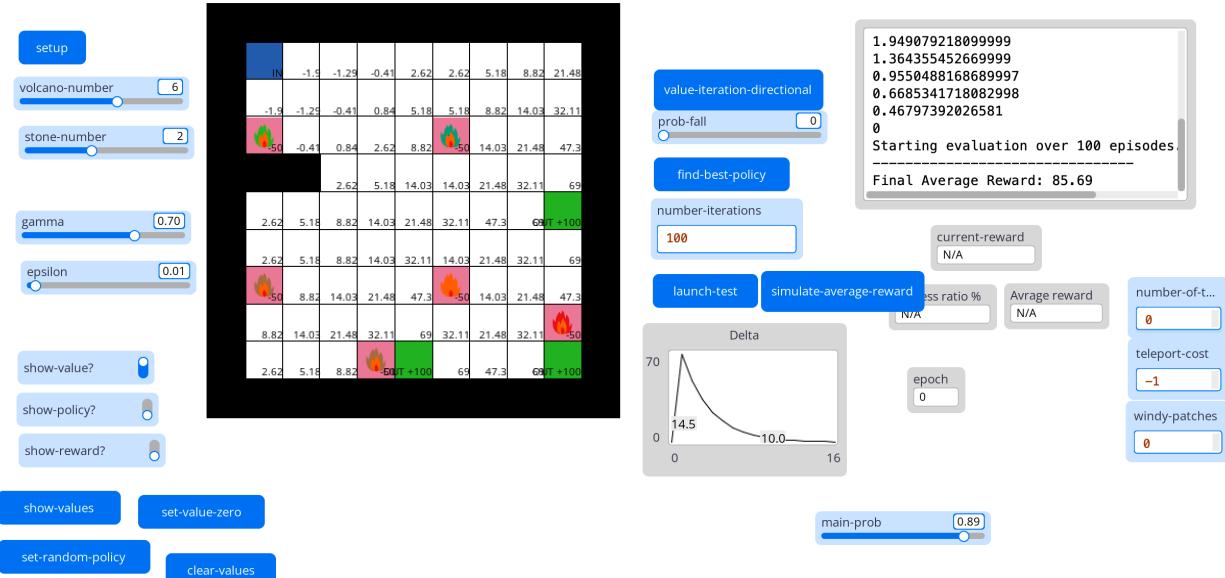


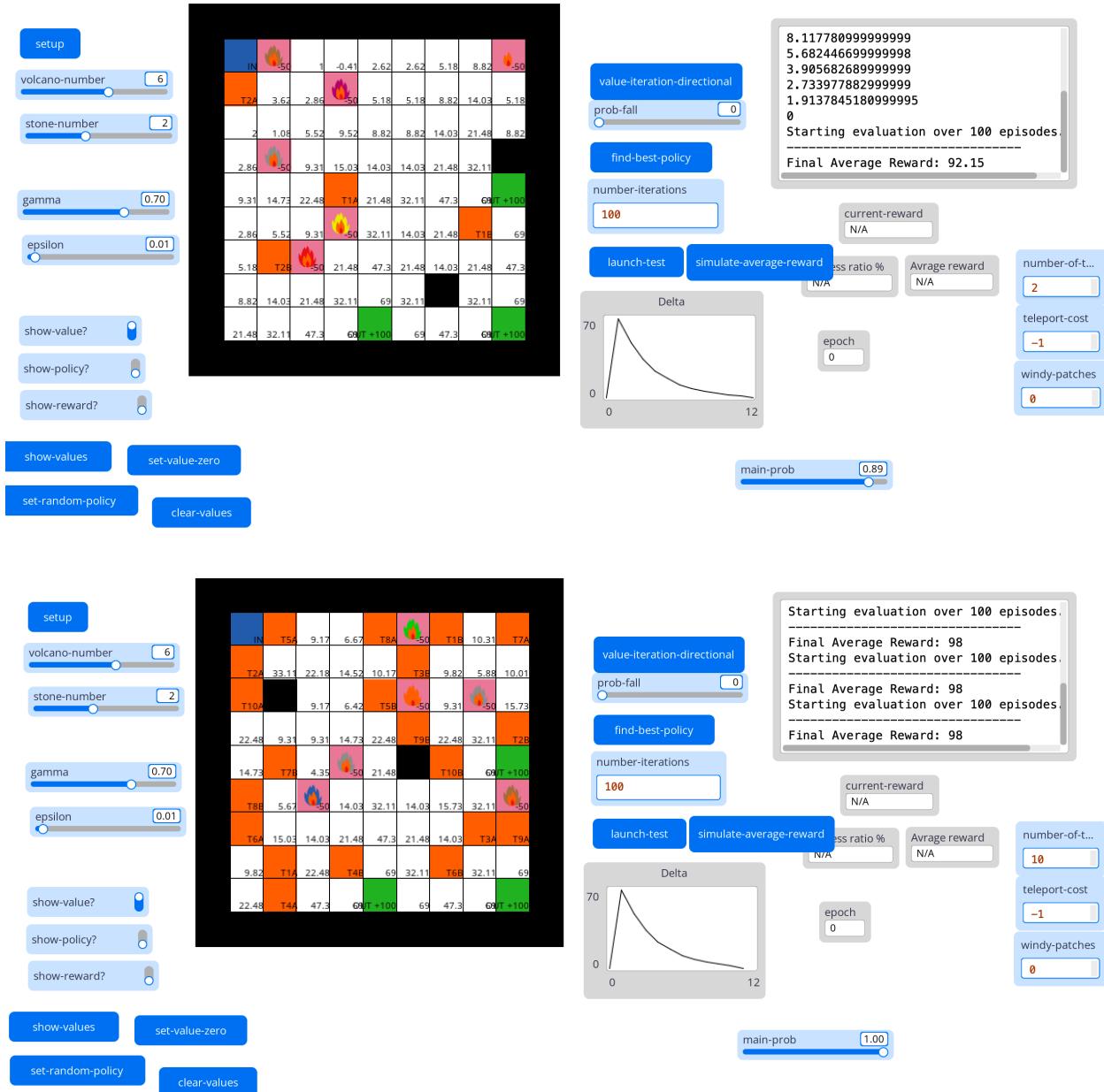




We can see that optimal learning rate is ≈ 0.5 . Obviously, the higher probability of success of given action the better final reward.

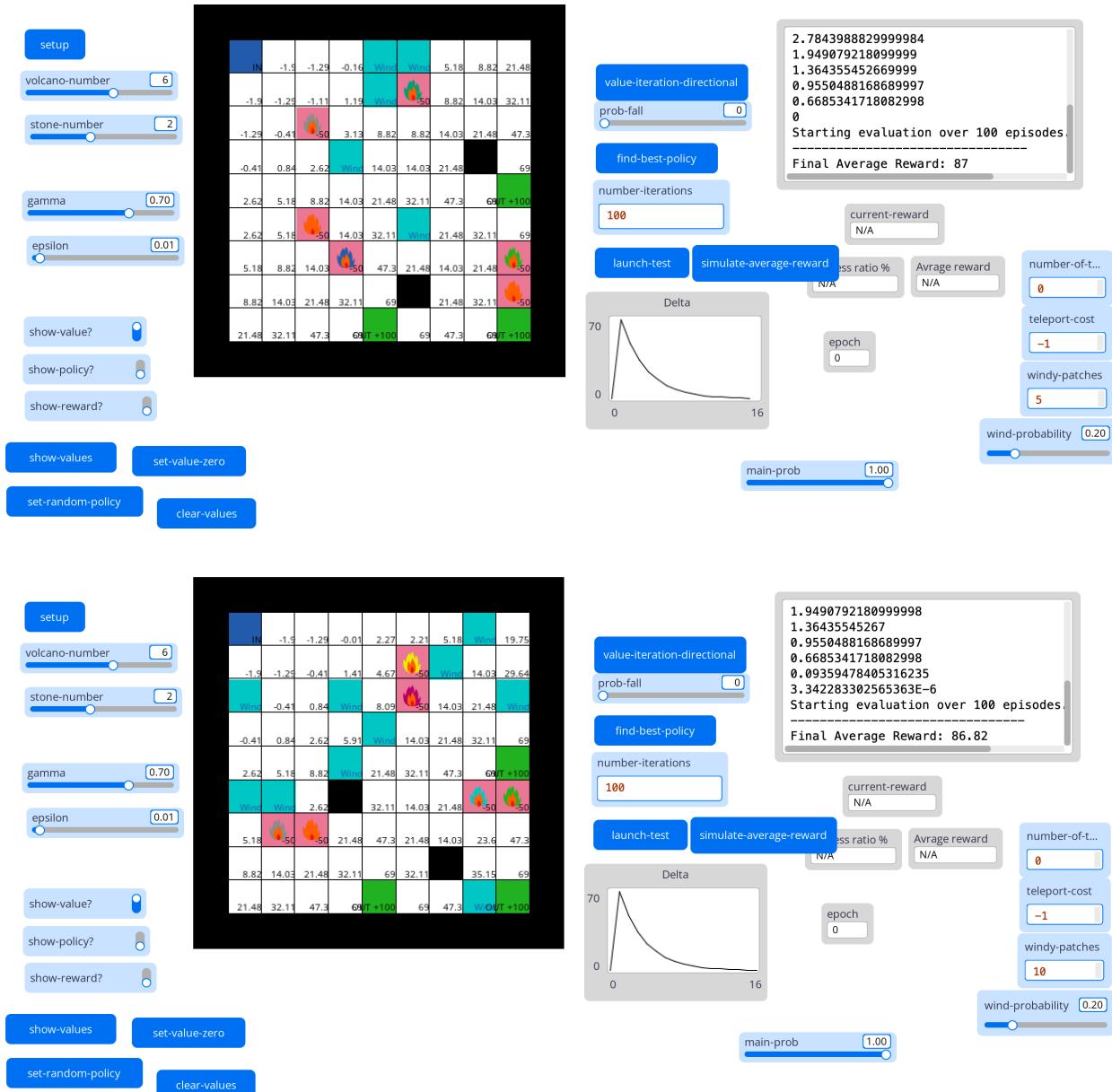
Teleports

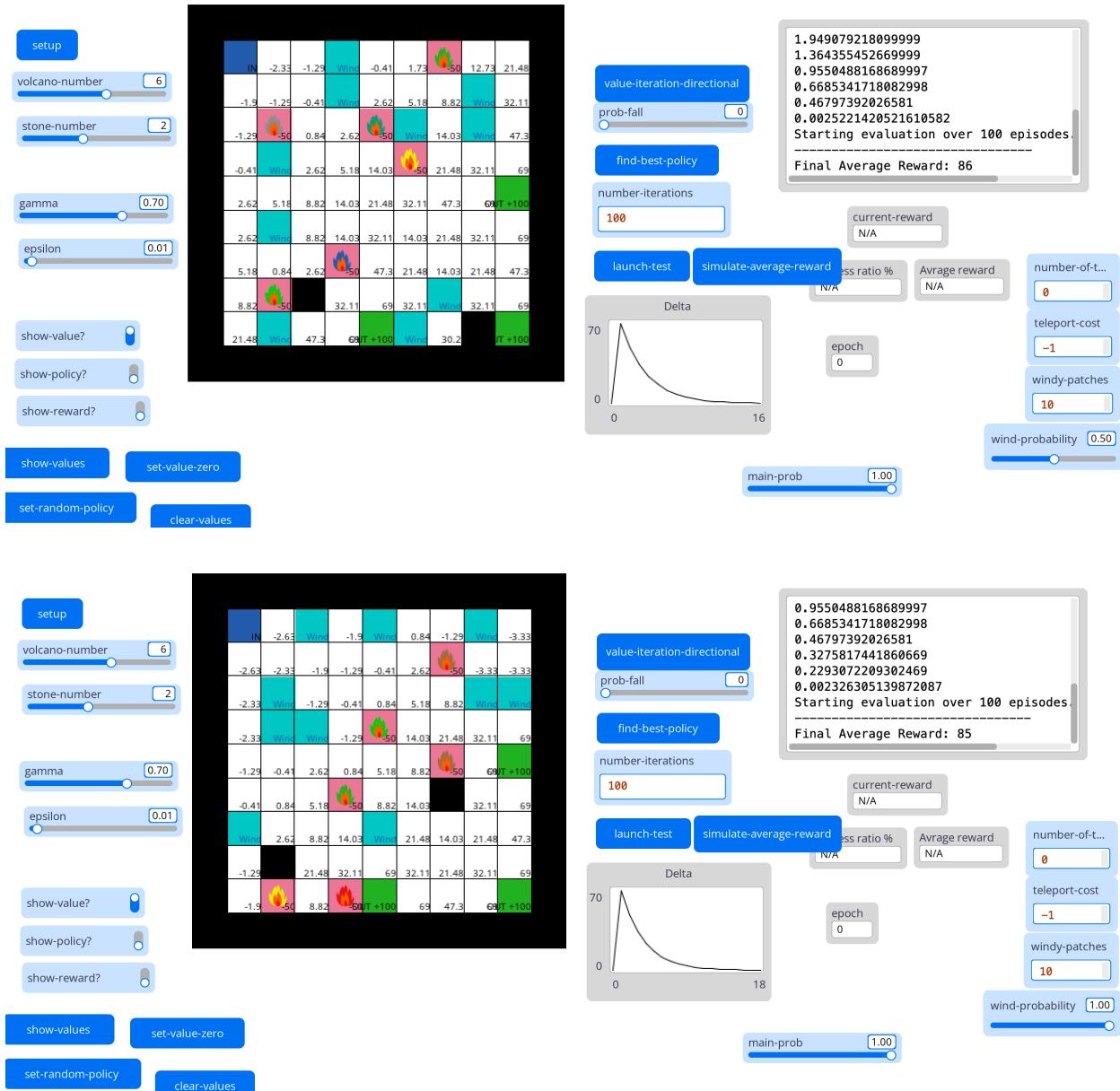


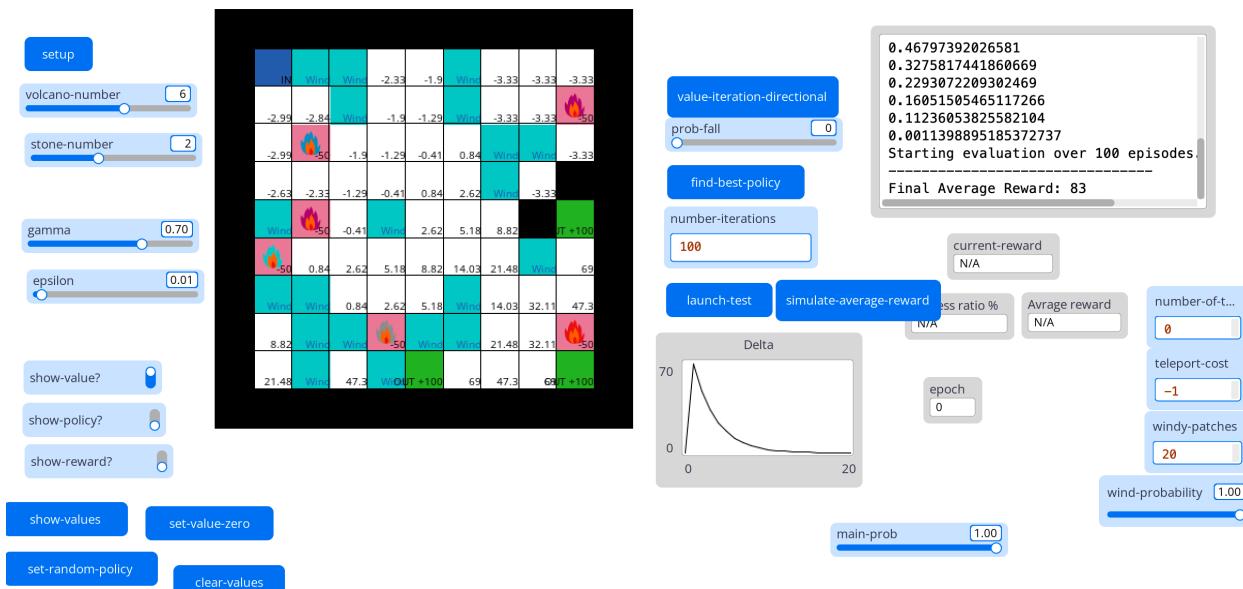
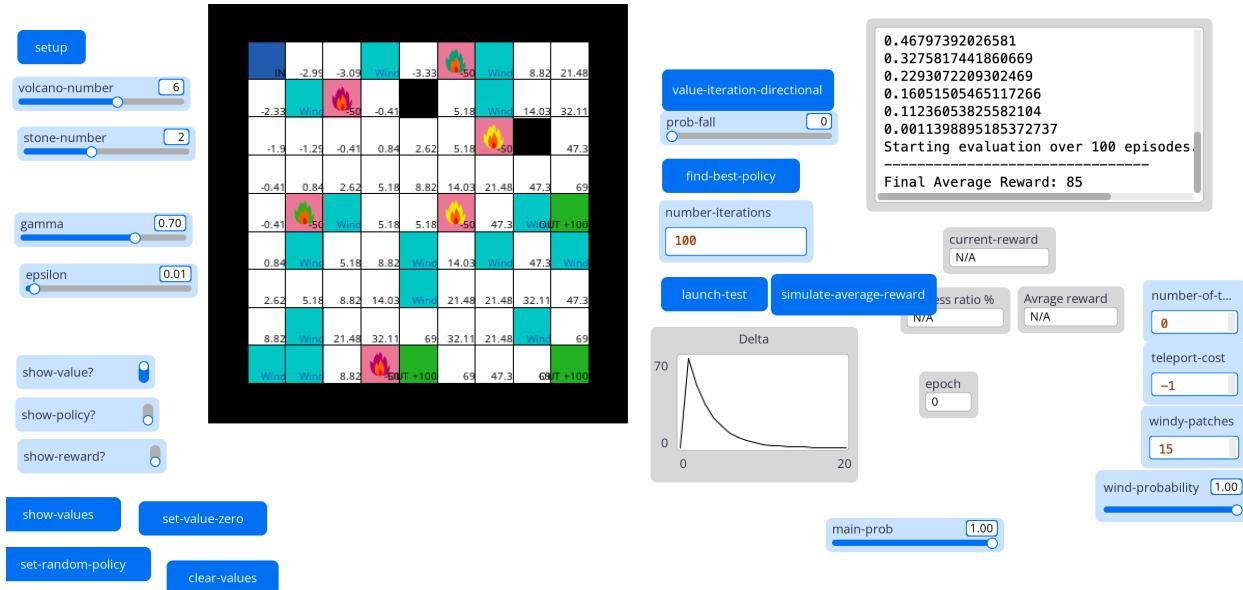


We can see, that average reward growth as the number of teleports growth. This is explained by the fact that agent learns position of teleport, which create a shortcut.

Wind





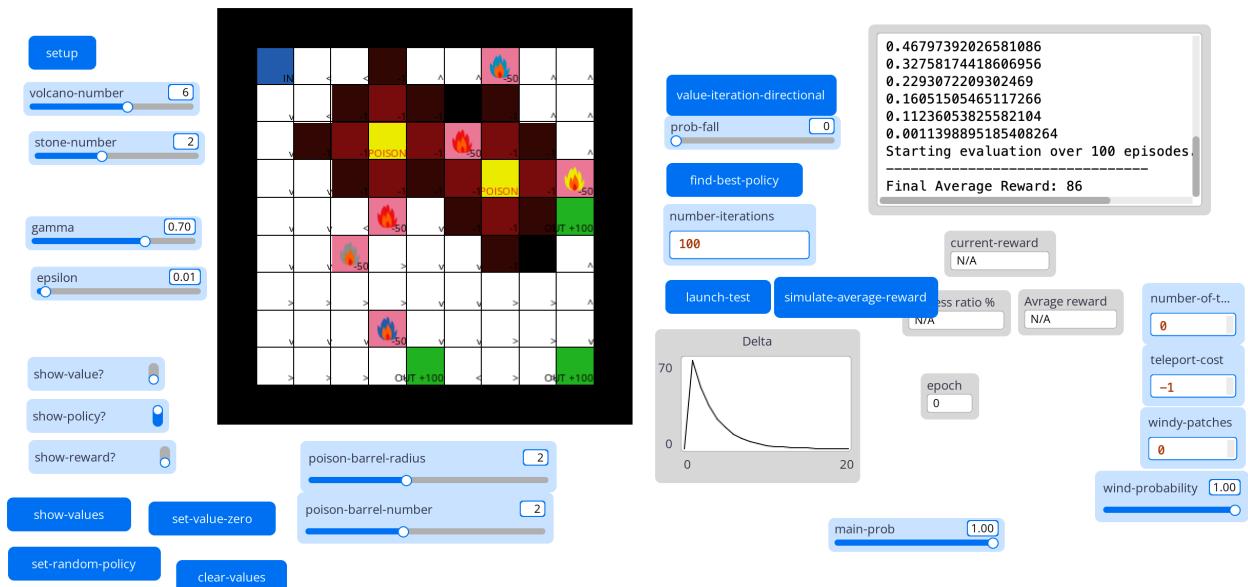
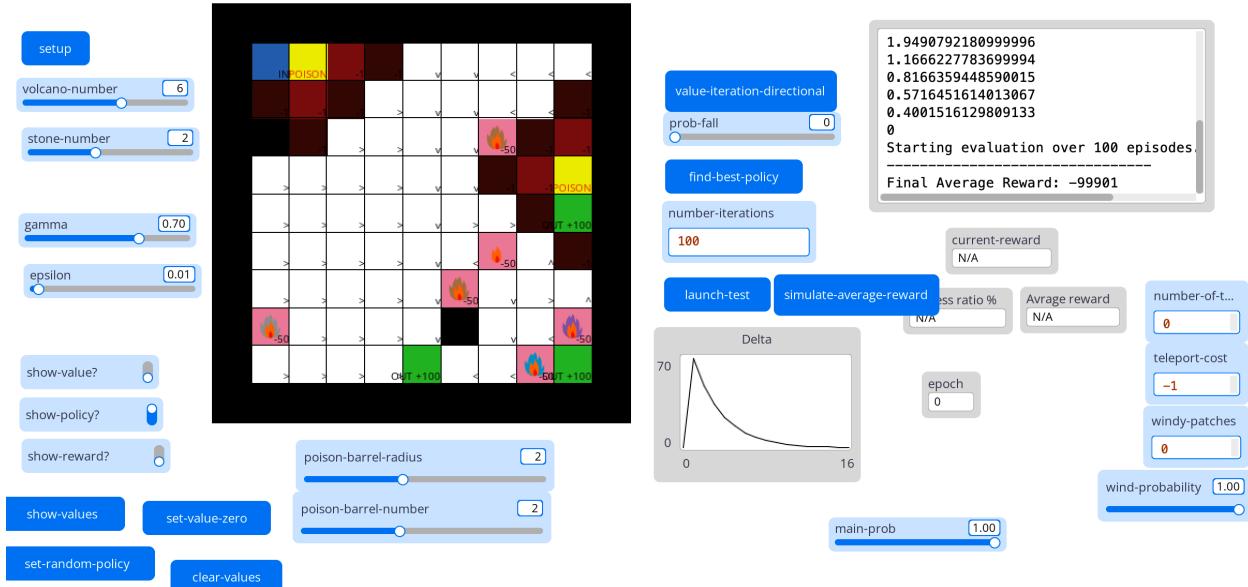


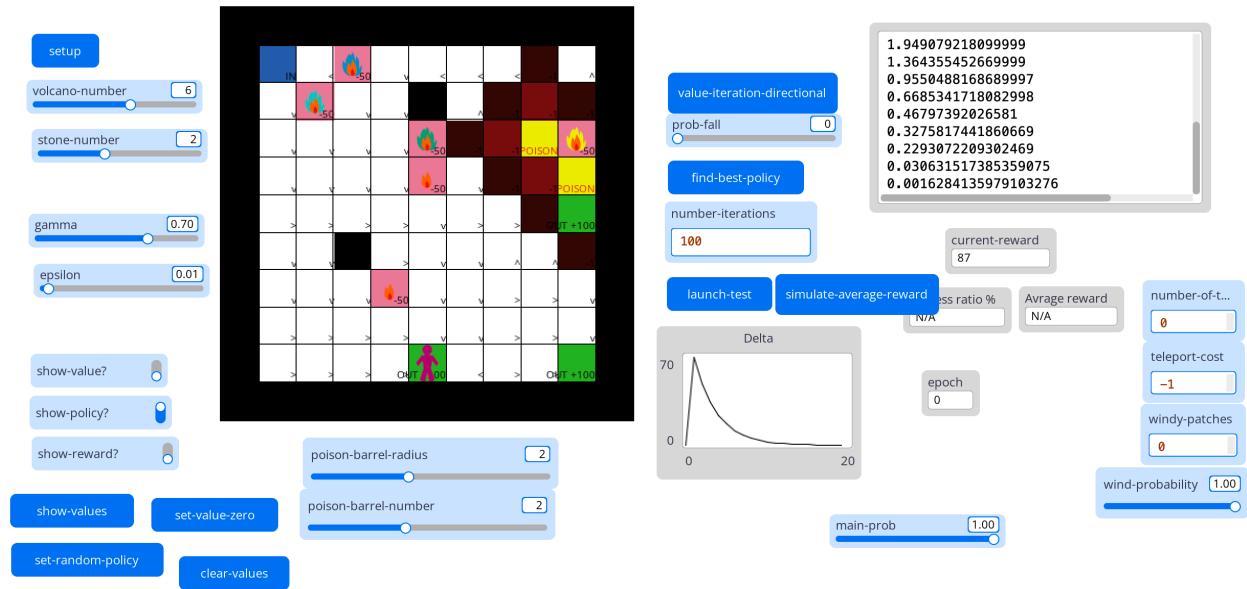
In general, I do not see that average reward decreases as number of windy blocks and probability increases. My assumption is that wind blows in the direction of OUT, thus agent learns how to escape "risky" windy blocks, and how to abuse a handful ones.

Poison

I implemented poison as described in the task. To include it in the VIA model, I add probability of being poisoned * punishment cost to expected probability.

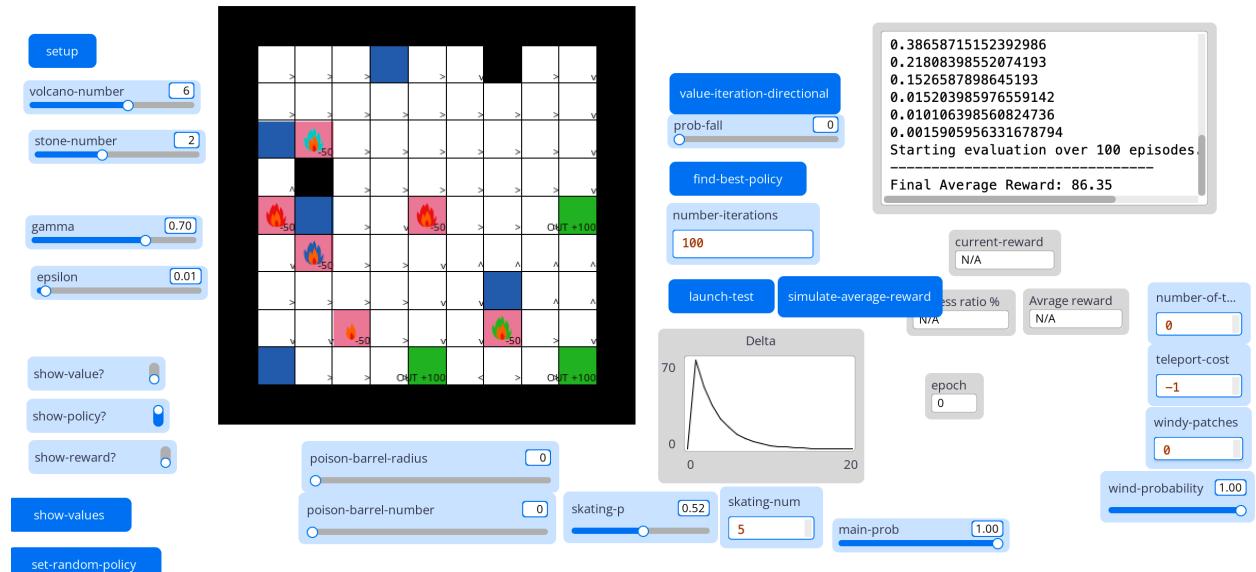
It is interesting that agent cannot escape starting point after solving policy, if starting point is surrounded by poison

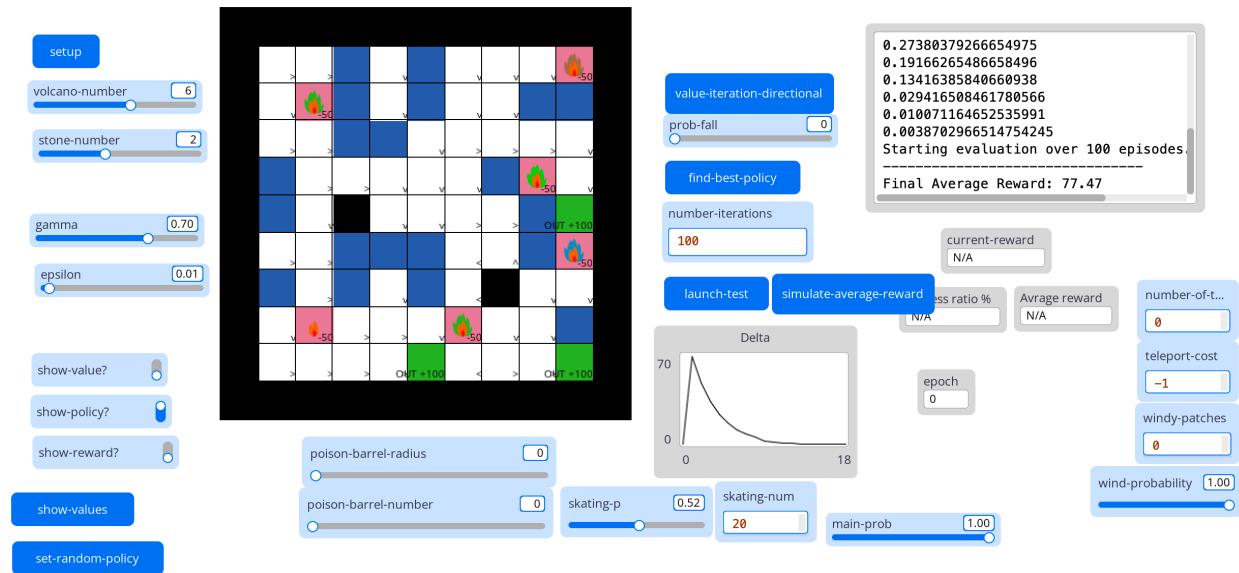
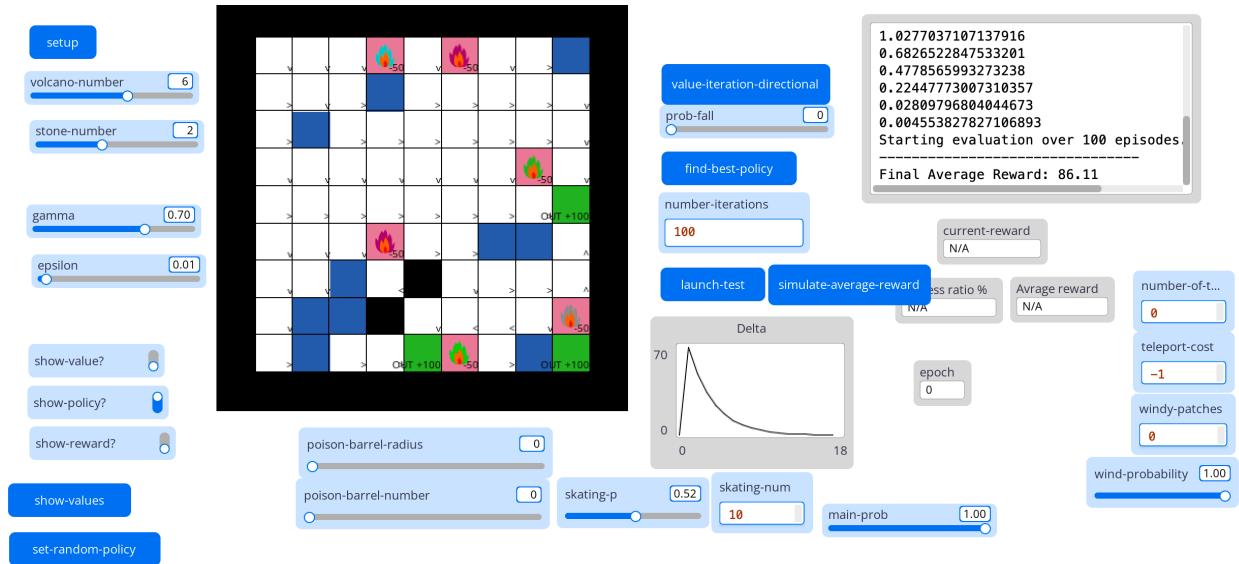


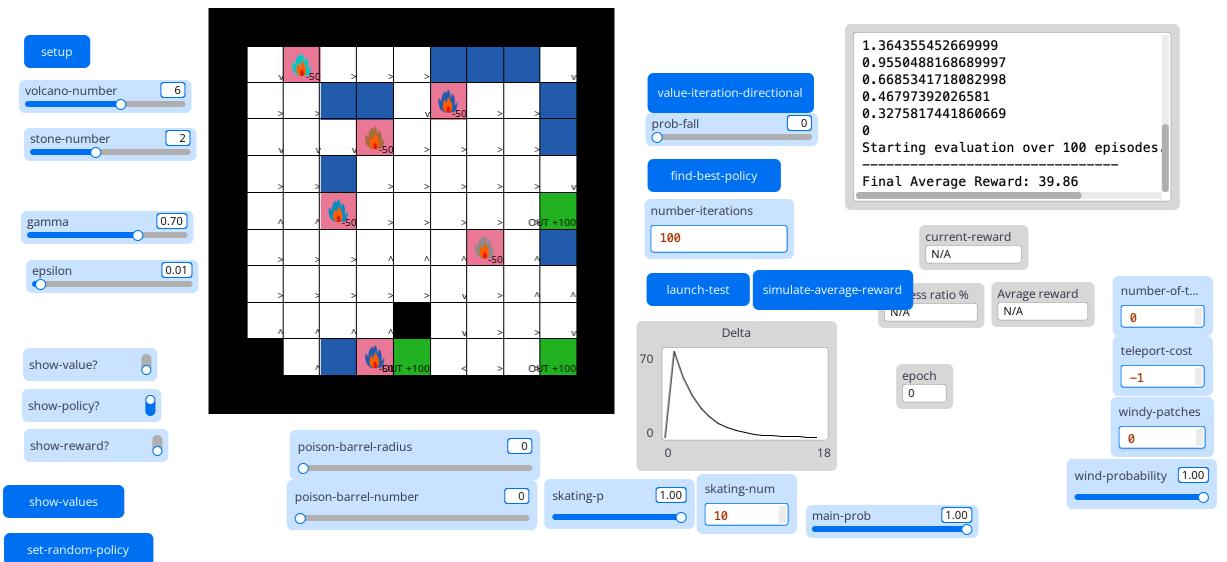
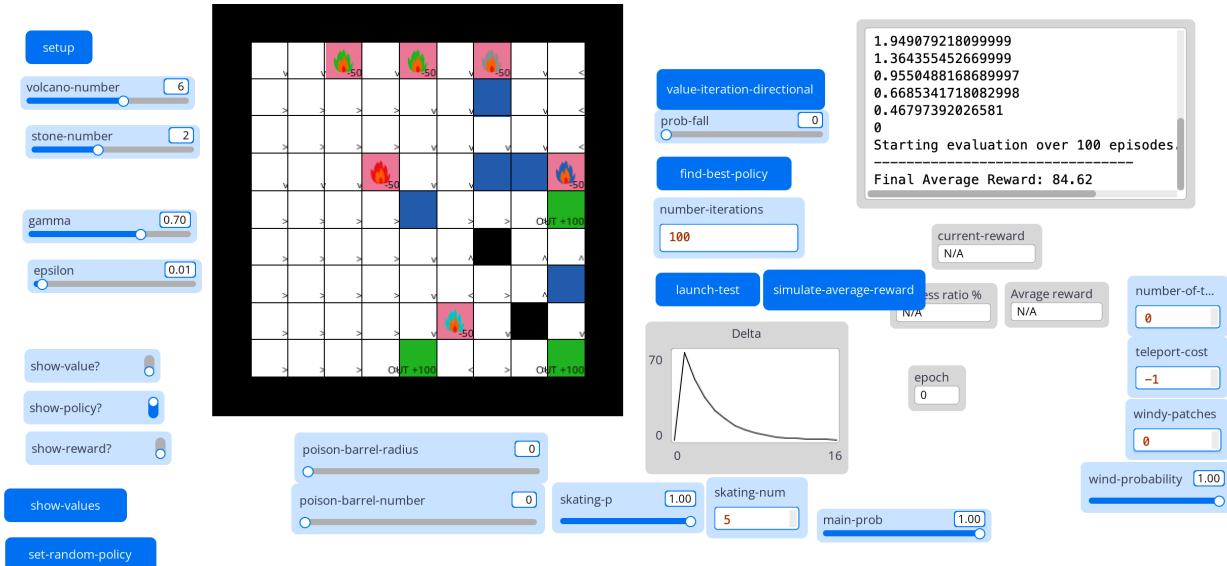


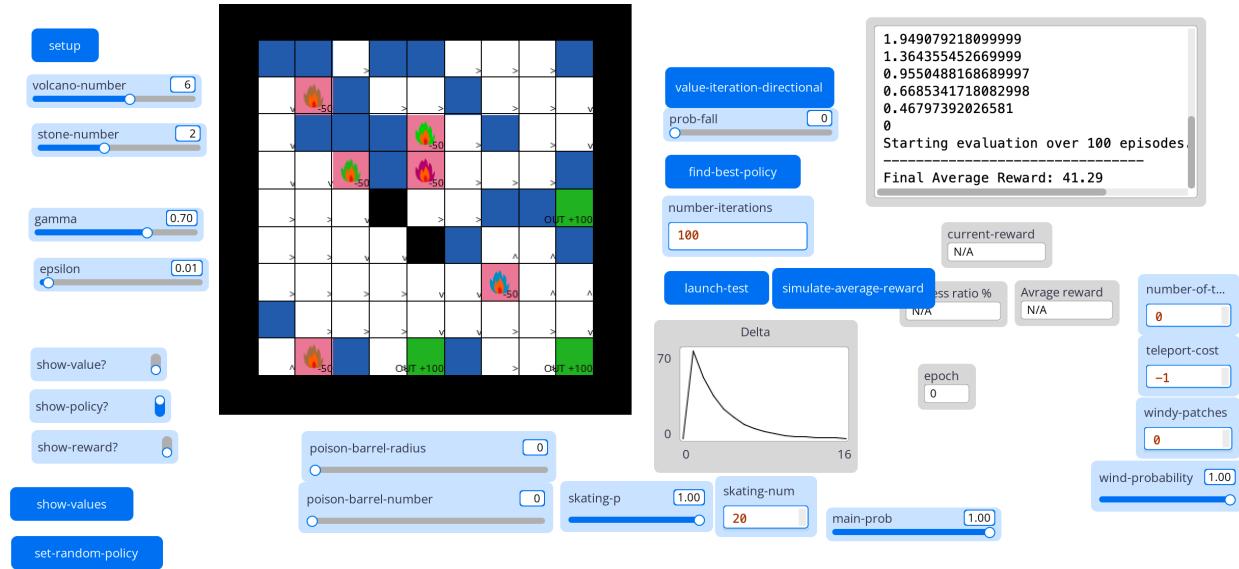
However, if starting point is not surrounded everything is fine.

Ice skating









Increasing number of skating patches or skating probability certainly seems to decrease average reward.

Additional task. Moving blocks.

