

Programming Assignment 3

Sai Praveen B

April 10, 2016

1 QA

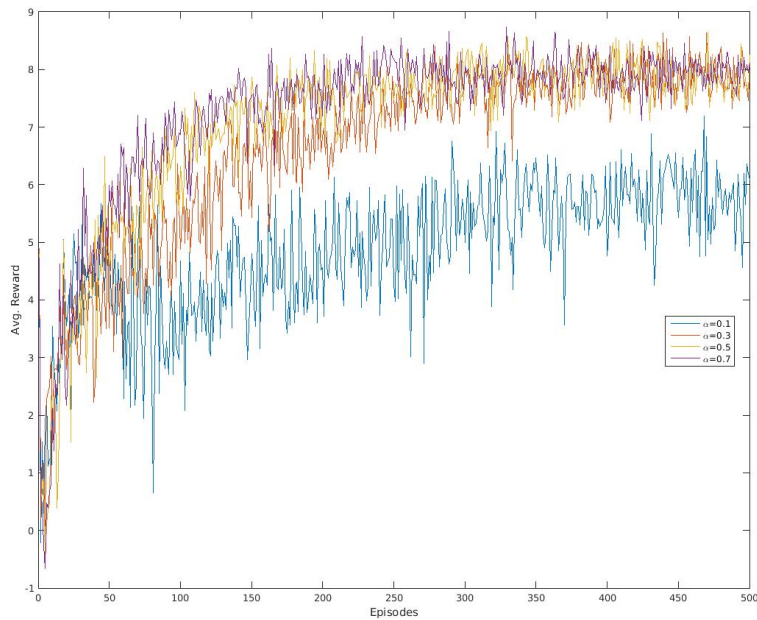
RL-glue for C++ has been used to simulate the agent and environment.

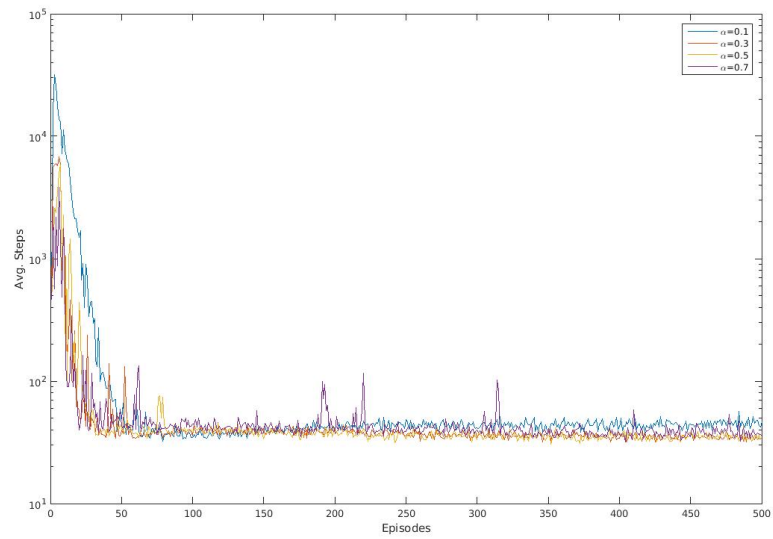
Instead of averaging over 50 runs, Averages are over 200 runs to reduce the noise(most of the methods, especially MC Policy gradient hav very high variance).

1.1 Q-learning

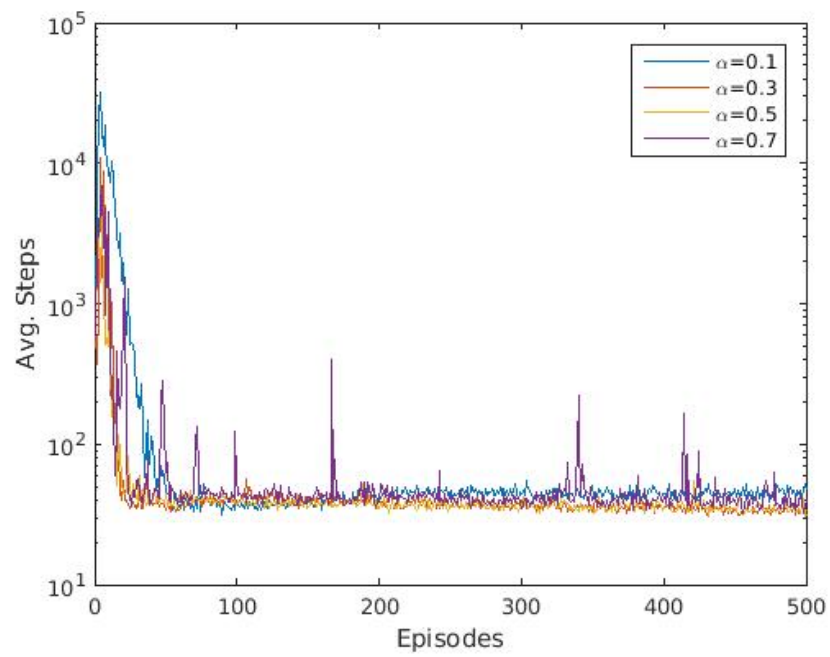
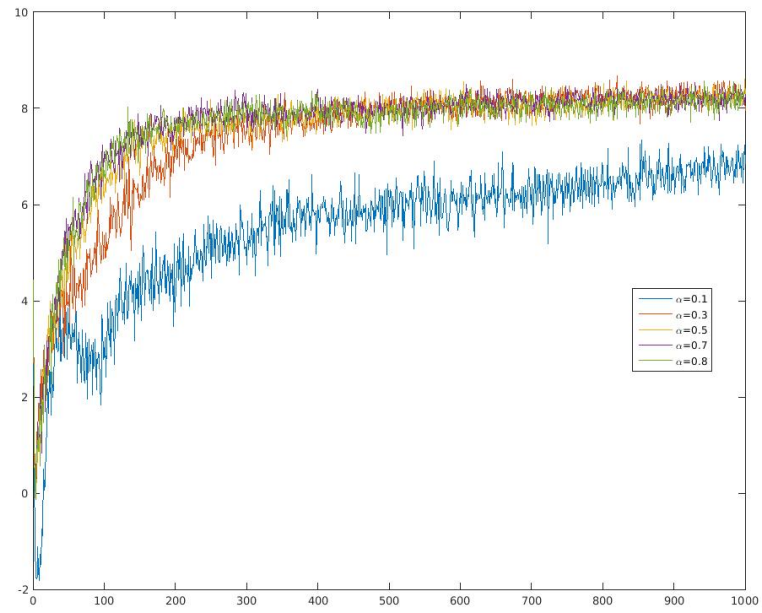
Q - learning has been implemented in C++(File: QAgent.cpp) The performance plots of the agent for each of the three targets are as follows:

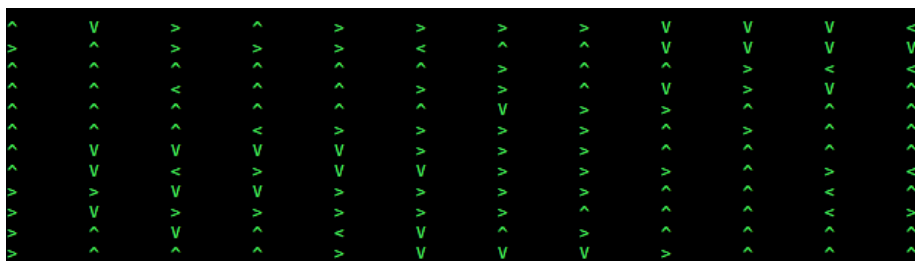
(Averages over 250 runs for various alpha) Target A:



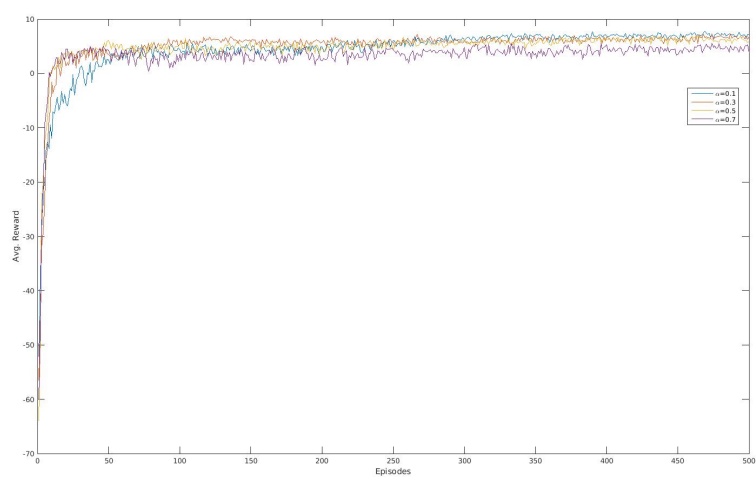


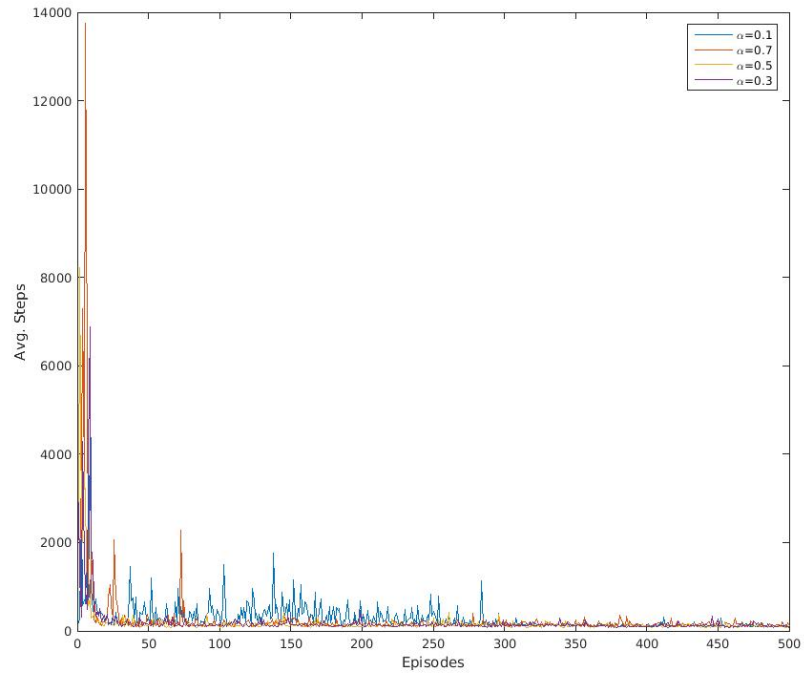
Target B:





Target C:

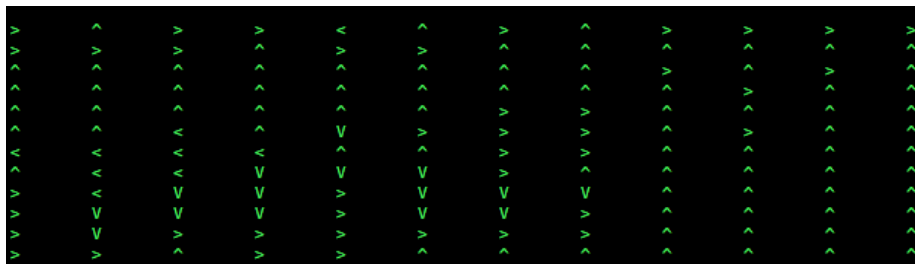
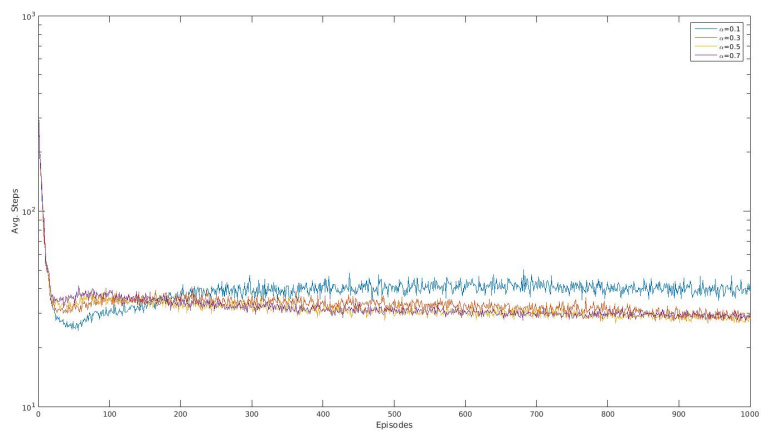
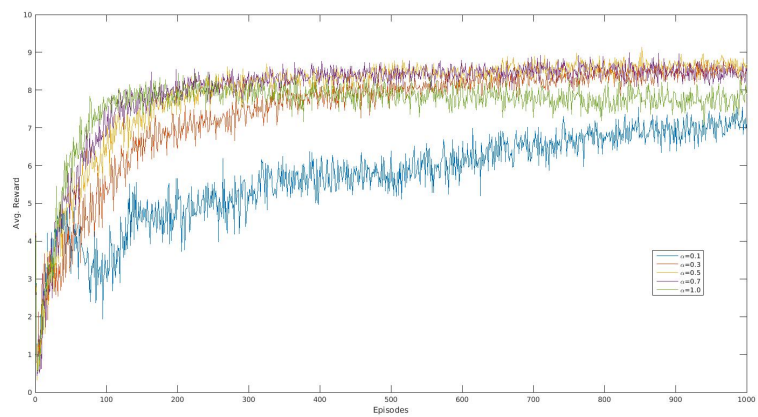




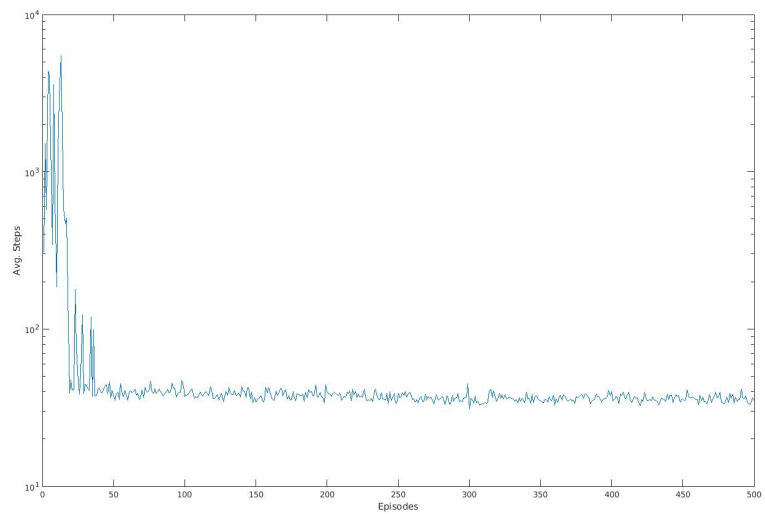
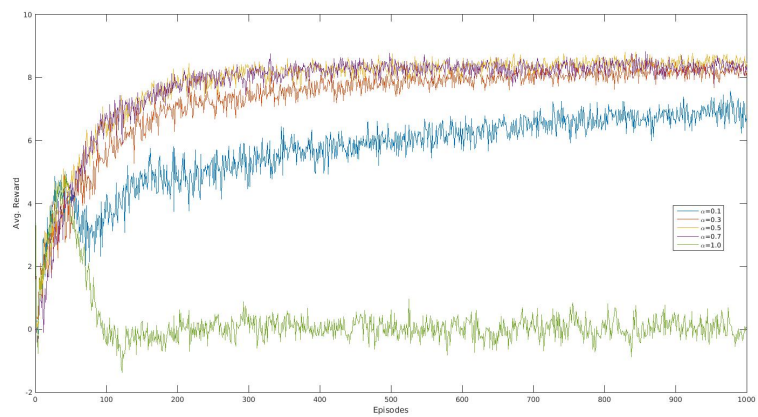
1.2 SARSA

The standard SARSA algorithm has been implemented in C++(File: SARSAAgent.cpp)

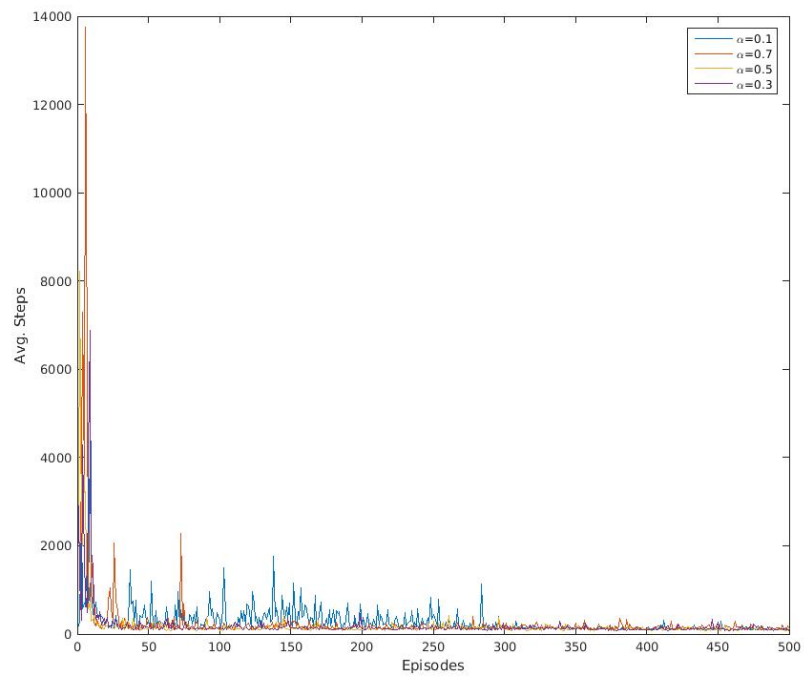
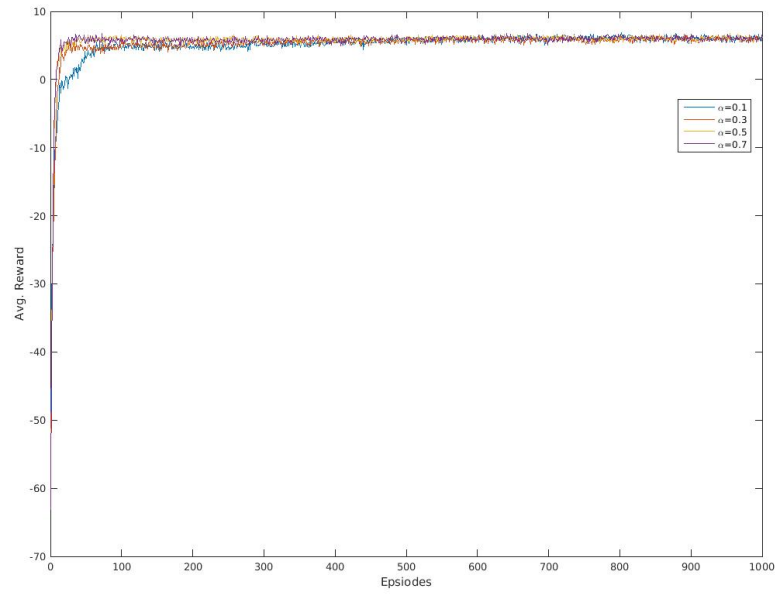
Target A:

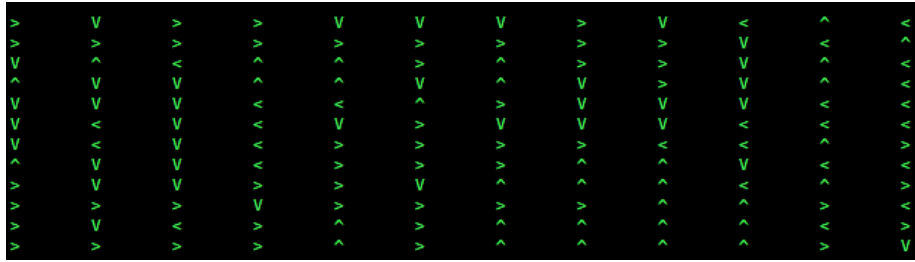


Target B:



Target C:





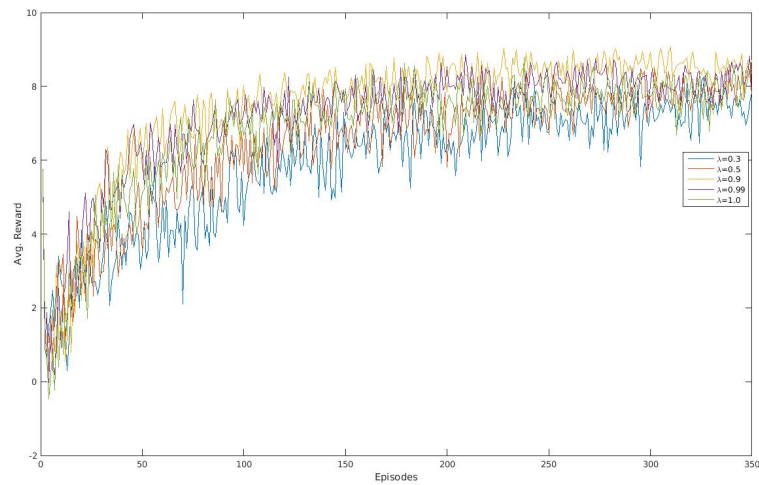
1.3 SARSA Lambda

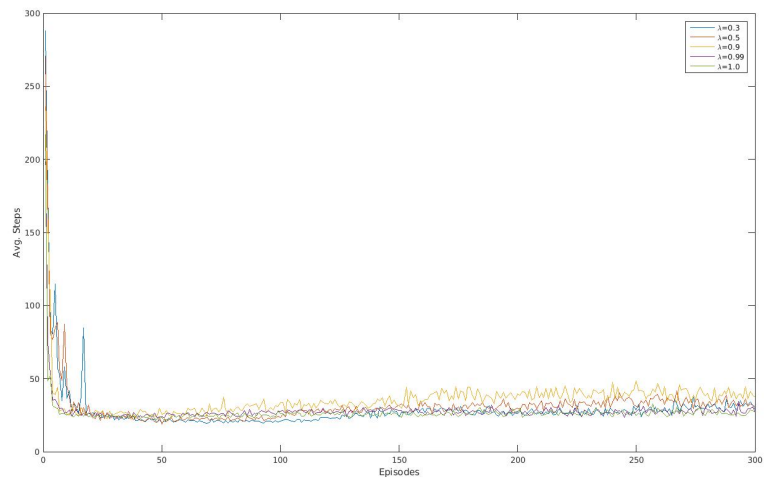
SARSA Lambda has been implemented with 2 methods. The first one uses a history list of episodes and thus implements Accumulating traces.(File: SarsaLambdaAgent.cpp) The second one uses a matrix of eligibilities and currently implements Replacing traces. (File: SarsaFSLambdaAgent.cpp)

The results shown below are averaged over 50 runs for Replacing Traces only.

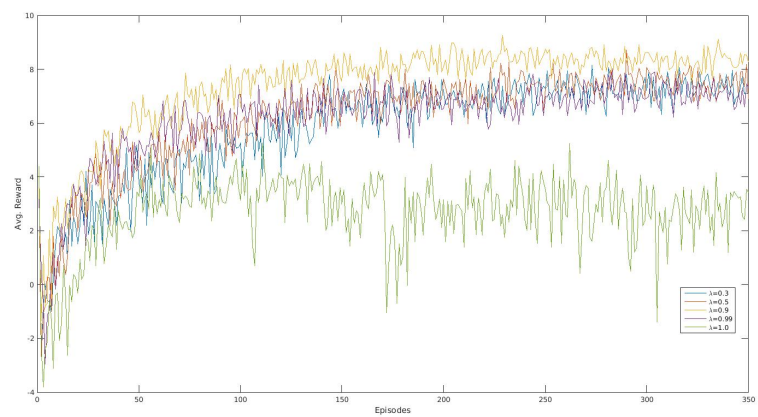
The results have been plotted for 5 values of $\lambda = 0.3, 0.5, 0.9, 0.99, 1.0$ Target

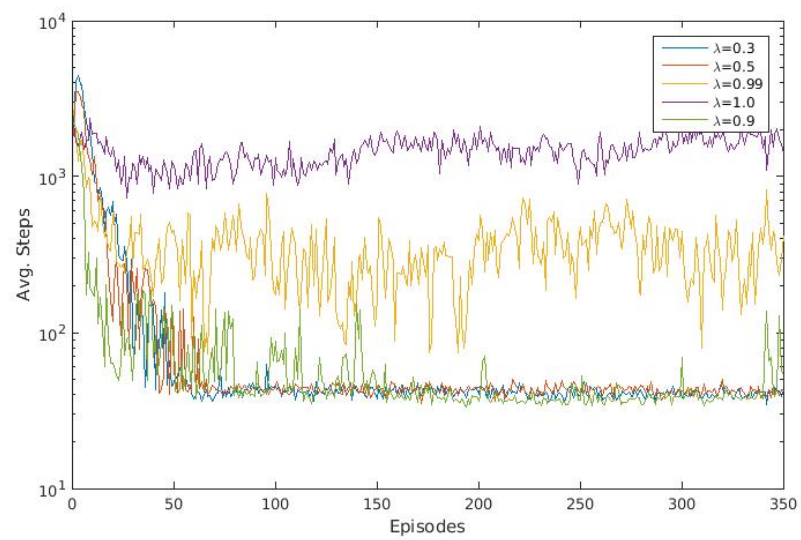
A:



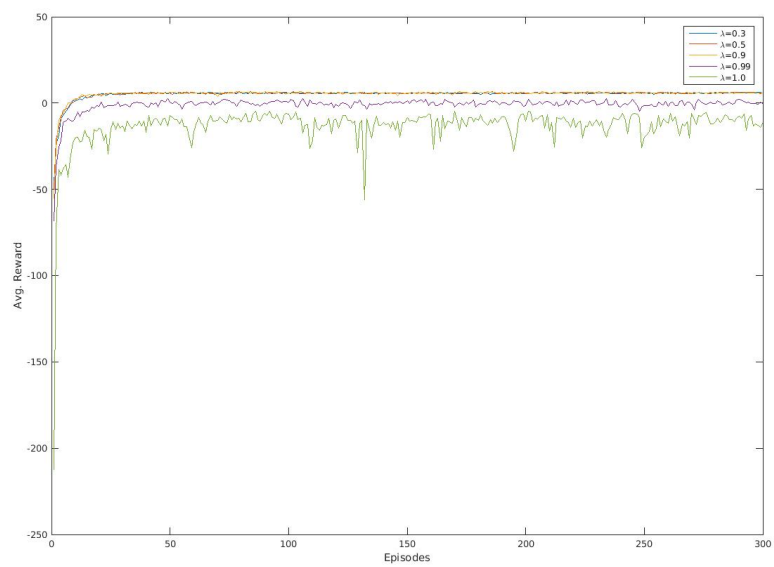


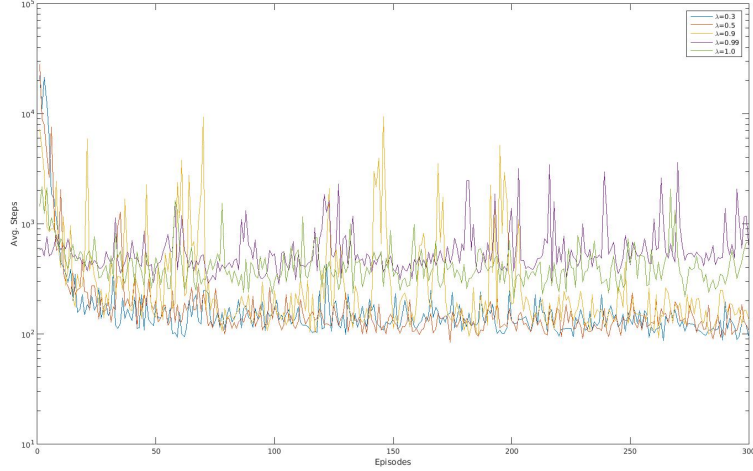
Target B:





Target C:





1.4 MC Policy Gradient

A Monte Carlo Policy gradient technique has been implemented in `MCGradientAgent.cpp`

We derive the update rules to the REINFORCE-like equations as follows:
Our function is:

$$\pi_{\theta}(s_t, a_t) = \frac{e^{\theta_{s_x, t, a_t} + \theta_{s_y, t, a_t}}}{\sum_b e^{\theta_{s_x, t, b} + \theta_{s_y, t, b}}}$$

Our target algorithm:

```
function REINFORCE
  Initialise  $\theta$  arbitrarily
  for each episode  $\{s_1, a_1, r_2, \dots, s_{T-1}, a_{T-1}, r_T\} \sim \pi_{\theta}$  do
    for  $t = 1$  to  $T - 1$  do
       $\theta \leftarrow \theta + \alpha \nabla_{\theta} \log \pi_{\theta}(s_t, a_t) v_t$ 
    end for
  end for
  return  $\theta$ 
end function
```

We now differentiate to get the gradient as:

$$\Delta_{\theta_{s_x, t, a_t}} \pi_{\theta}(s'_t, b_t) = 1 - \pi_{\theta}(s_t, a_t), s_{x, t} = s'_{x, t}, a_t = b_t$$

$$\Delta_{\theta_{s_{x,t},a_t}} \pi_{\theta}(s'_t, b_t) = 1 - \pi_{\theta}(s_t, a_t), s_{x,t} = s'_{x,t}, a_t = b_t$$

$$\Delta_{\theta_{s_{x,t},a_t}} \pi_{\theta}(s'_t, b_t) = 0, otherwise$$

The same rules apply to y as well.

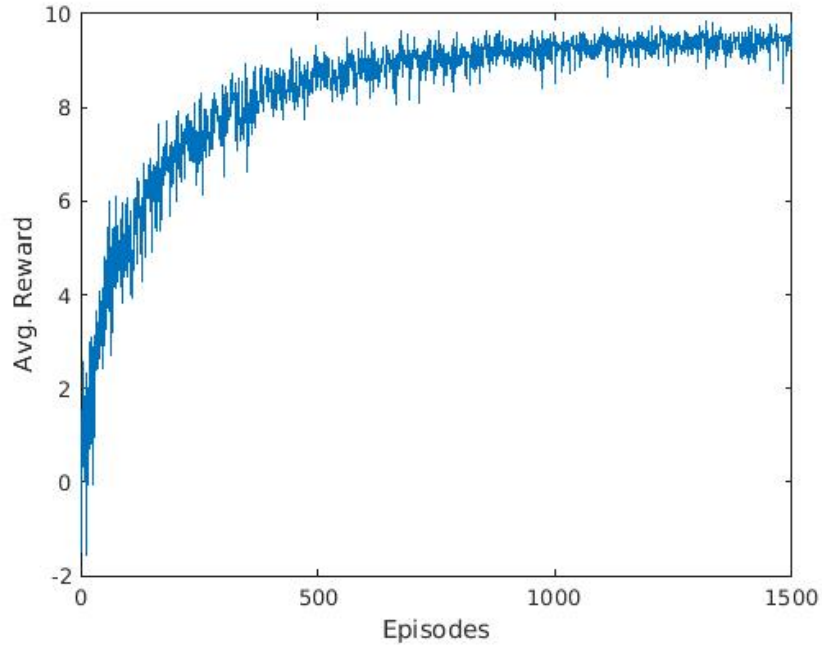
The gradient is factored with an α (learning rate) which has been fixed at $\alpha = 0.012$ after much trail and error. The method uses a per row/column per action parametrization which is inherently bad for the general case.

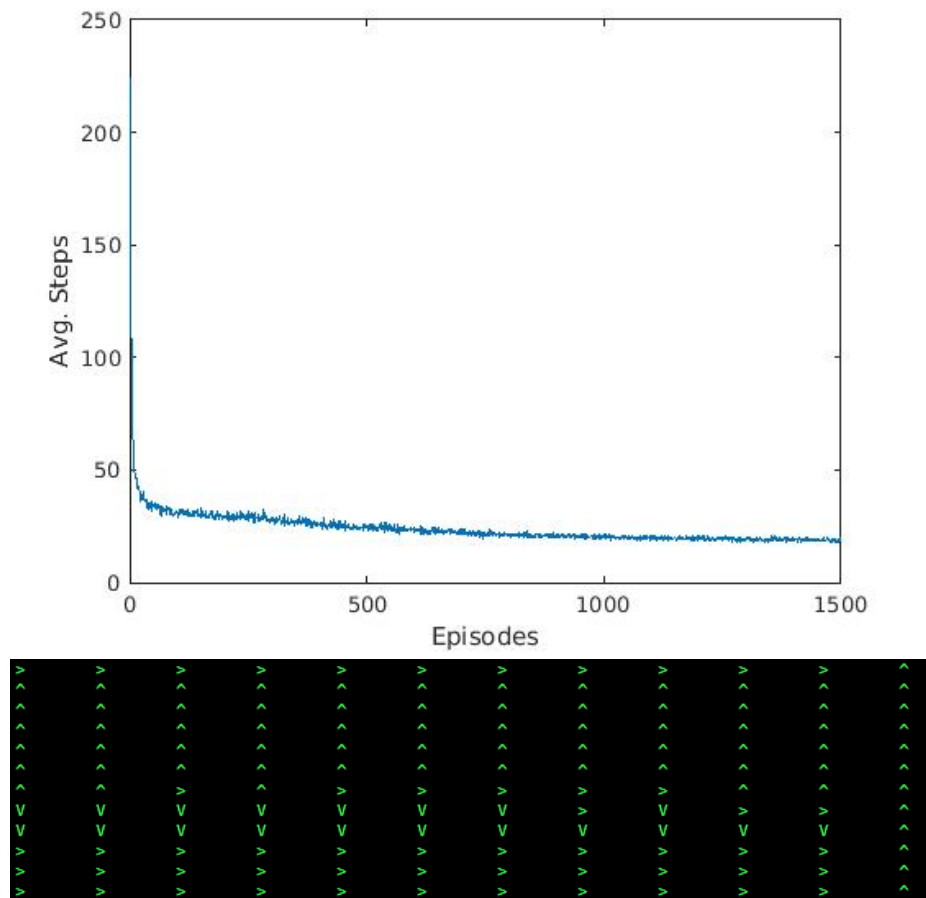
Target A is very easy and has an optimal solution that can be represented using the parametrization given. It gets an average reward of around 9. Avg. Step count of around 20.

Target B is considerably tougher given that the row has other negative cells along it and an attempt to optimize near the target can lead to bad decisions elsewhere. It manages an average reward of around 7 and an average step count of around 50.

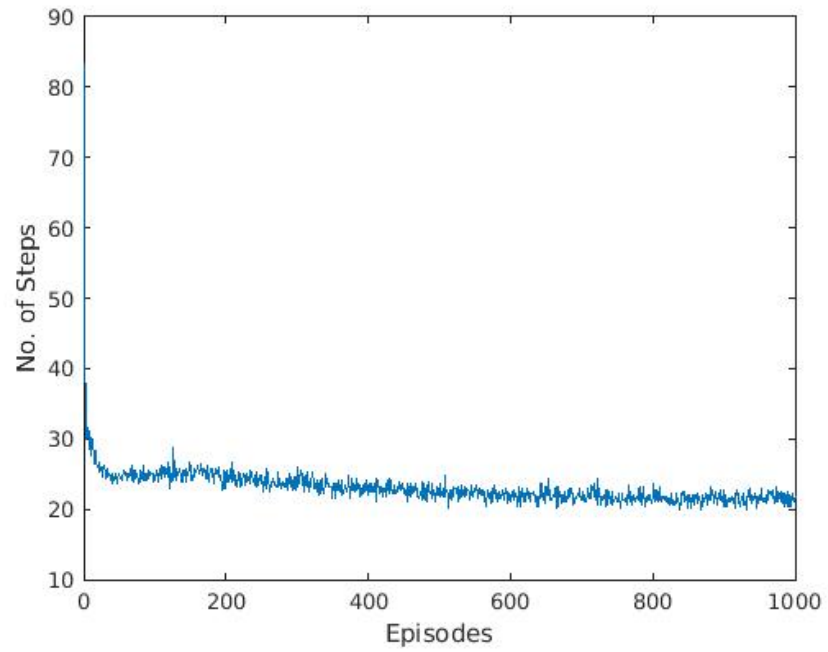
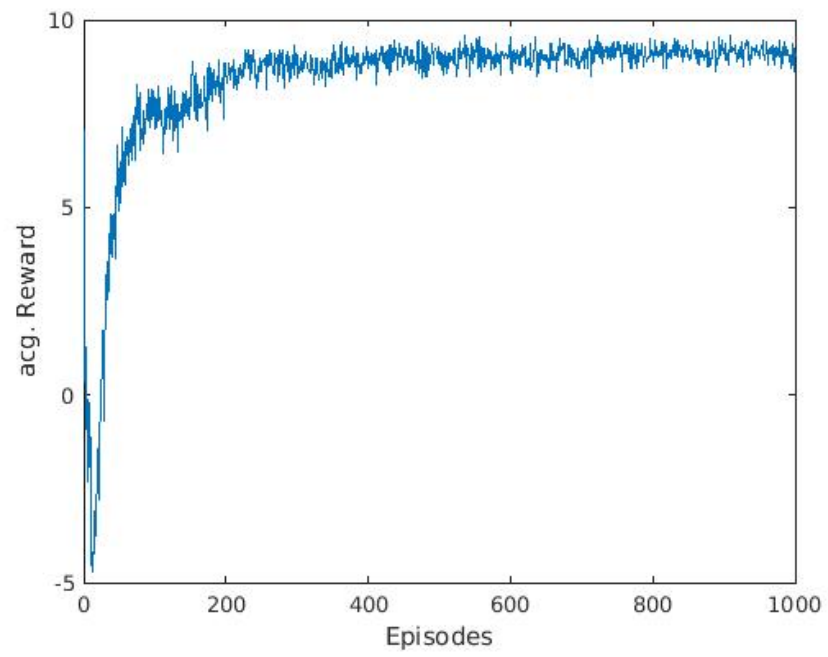
Target C is very difficult as it's in the middle of the puddle with high negatives in both the row and column. This leads to a variety of issues including the agent often getting stuck due to bad loops formed due to ripples that occur because of the non-local nature of the parametrization. It only manages a reward of around 4 and an even worse step count of around 400.

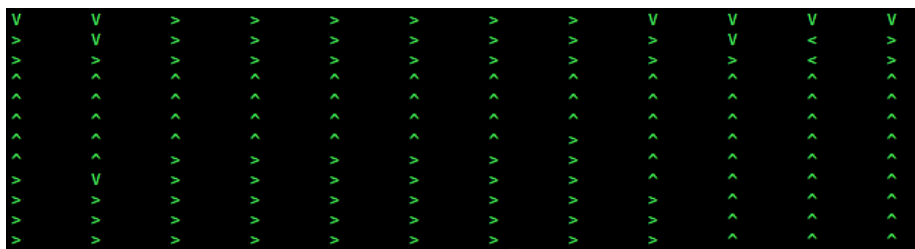
Target A:



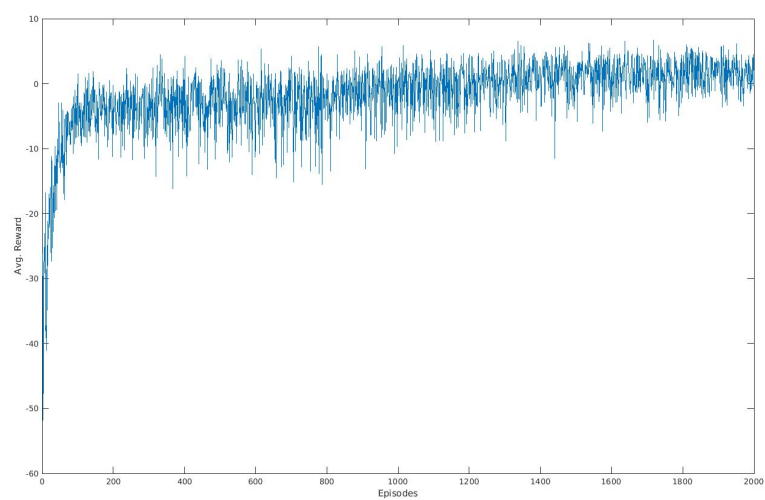


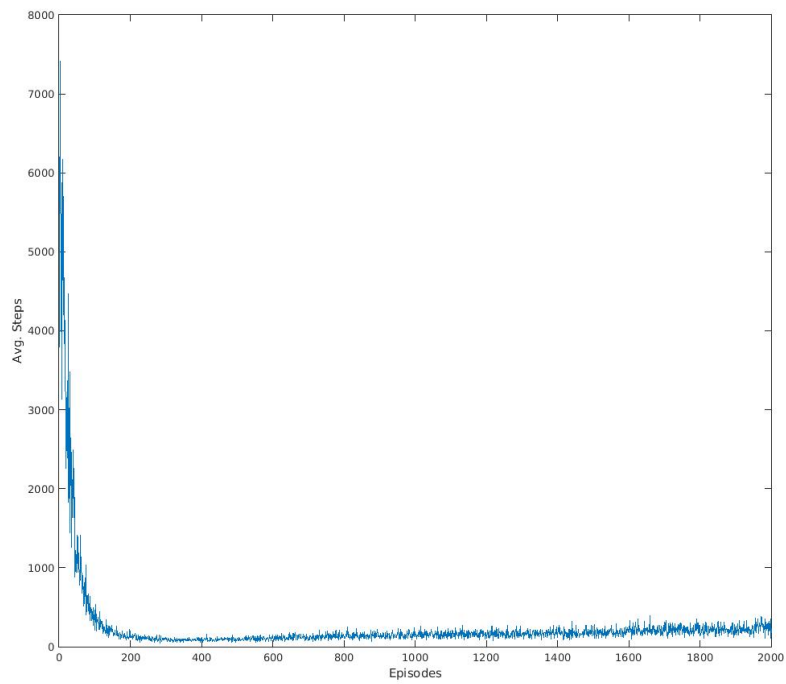
Target B:





Target C:





The parametrization is quite bad as it oscillates in the case of C sometimes and sometimes causes the agent to get stuck in a loop.