



In The Name Of God

HW06

Advanced Neuroscience

MohammadAmin Alamalhoda
97102099

■ Part1 - Model Free RL

□ Explanation on the used Method

I have used the **Model-Free** method for implementing the RL algorithm. The used equations for updating the states values are written down in the following:

$$lr_t = \frac{0.1}{1+\log t}$$
$$\delta_t = r_t + \gamma V(S_{t+1}) - V(S_t)$$

Where r_t is:

$$r_t(x, y) = \begin{cases} 2 & \text{if reward in x,y} \\ -2 & \text{if punishment in x,y} \\ 0 & \text{otherwise} \end{cases}$$

and

$$V(S_t) = V(S_t) + \delta_t lr_t$$
$$m_i = m_i + \delta_t (\delta_i - P(a_i))$$

Where δ_i and $P(a_i)$ are:

$$\delta_i = \begin{cases} 1 & \text{if } a_i \text{ is chosen} \\ 0 & \text{otherwise} \end{cases}$$
$$P(a_i) = \frac{e^{m_i}}{\sum_{j=1}^4 e^{m_j}}$$

□ Q01 and Q02

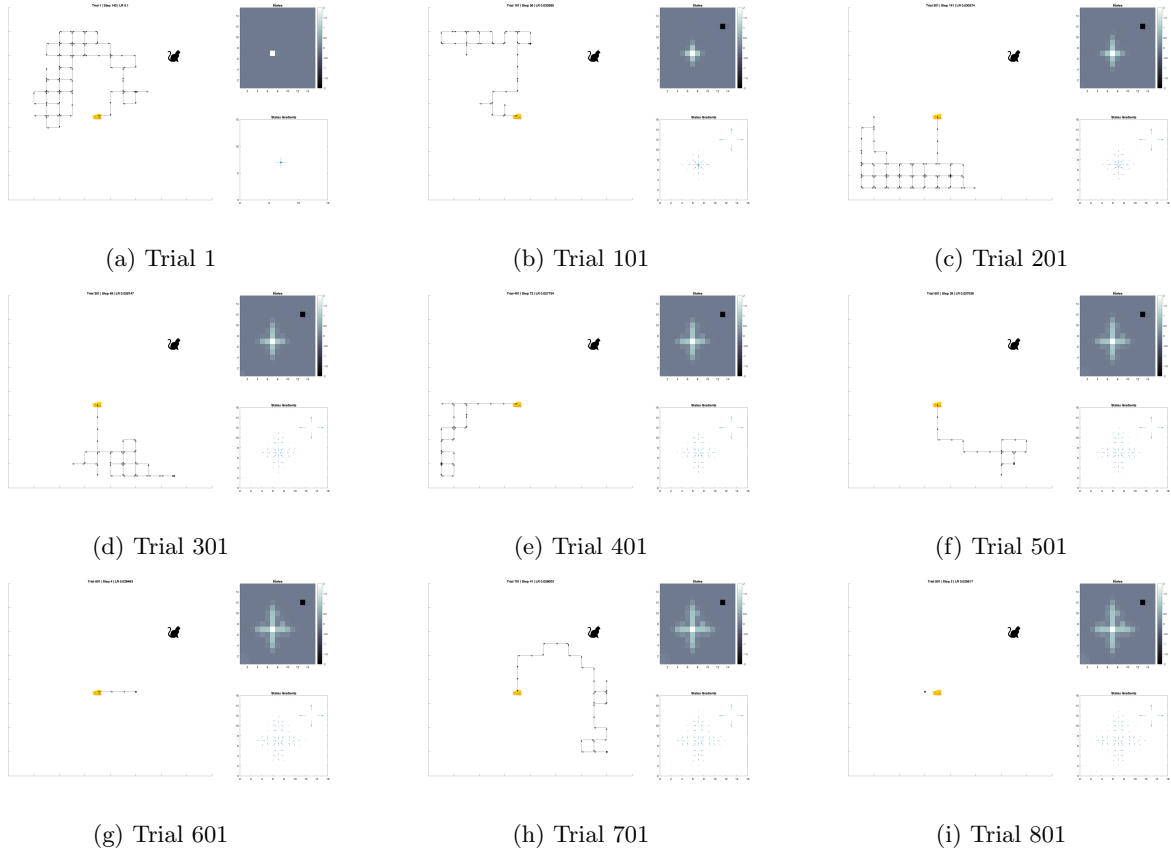


Figure 1: Agent passed way, States values and States gradients in some of the trials. $LR = \frac{0.1}{1+\log t}$, $\gamma = 1$

As can be seen in Figure 1, the agents learn the value of the states over the trials. It is noteworthy to mention that because my method in implementing the RL is Model-Free, the agent needs more trials to learn the value of the states and the transition probabilities.

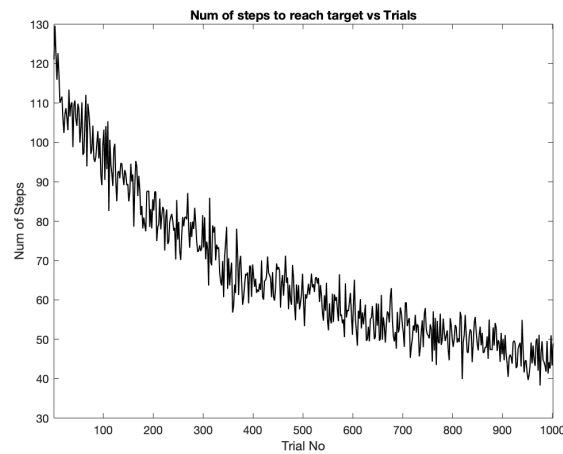


Figure 2: Number of steps to reach the target vs number of trials of learning. I run the model 200 times and average over all of the vectors of number of the steps in each of the iterations.

Figure 2 shows the relation between number of the training trials and the passed steps to reach the target. As can be seen, the passed steps to reach the target decreases over the trials.

□ Q03

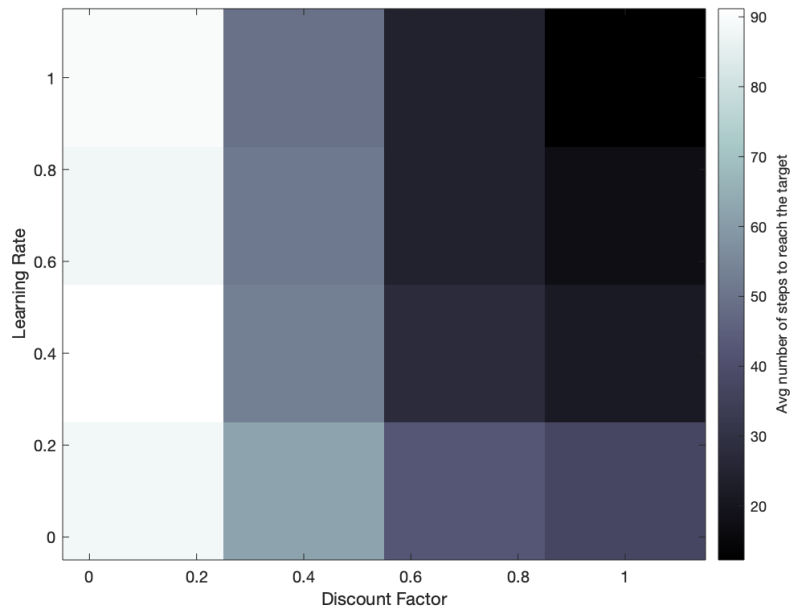


Figure 3: Average of passed steps to reach one target of the last 50 trials for different learning rates and discounting factors

□ Q04

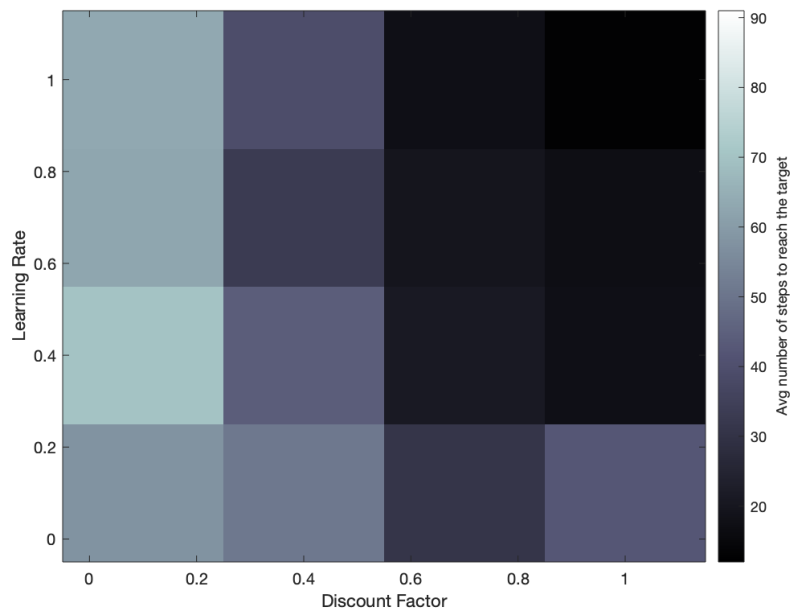


Figure 4: Average of passed steps to reach one of the two targets of the last 50 trials for different learning rates and discounting factors

As can be seen in Figures 3 and 4, bigger learning rates or discount factors tend to decrement in the number of passed steps to reach the target. Also, by adding another target the chance of the agent to reach a target increases and as can be seen in the Figure 4, the overall number of passed steps to reach one of the targets have been decreased in comparison with Figure 3.



■ Part2 - Model Free RL using TD rule

□ Explanation on the used Method

I have used the **Model-Free** method of TD rule for implementing the RL algorithm. The used equations for updating the states values are written down in the following:

$$lr_t = \frac{0.1}{1+\log t}$$
$$\delta_t = r_t + \gamma V(S_{t+1}) - V(S_t)$$

Where r_t is:

$$r_t(x, y) = \begin{cases} 2 & \text{if reward in x,y} \\ -2 & \text{if punishment in x,y} \\ 0 & \text{otherwise} \end{cases}$$

and

$$V(S_{t-t'}) = V(S_{t-t'}) + \lambda^{t'} \delta_t lr_t, \forall t' = 1, 2, \dots$$
$$m_{i(S_{t'})} = m_{i(S_{t'})} + \lambda^{t'} \delta_t (\delta_i - P(a_i))$$

Where δ_i and $P(a_i)$ are:

$$\delta_i = \begin{cases} 1 & \text{if } a_i \text{ is chosen} \\ 0 & \text{otherwise} \end{cases}$$
$$P(a_i) = \frac{e^{m_i}}{\sum_{j=1}^4 e^{m_j}}$$

□ Q05

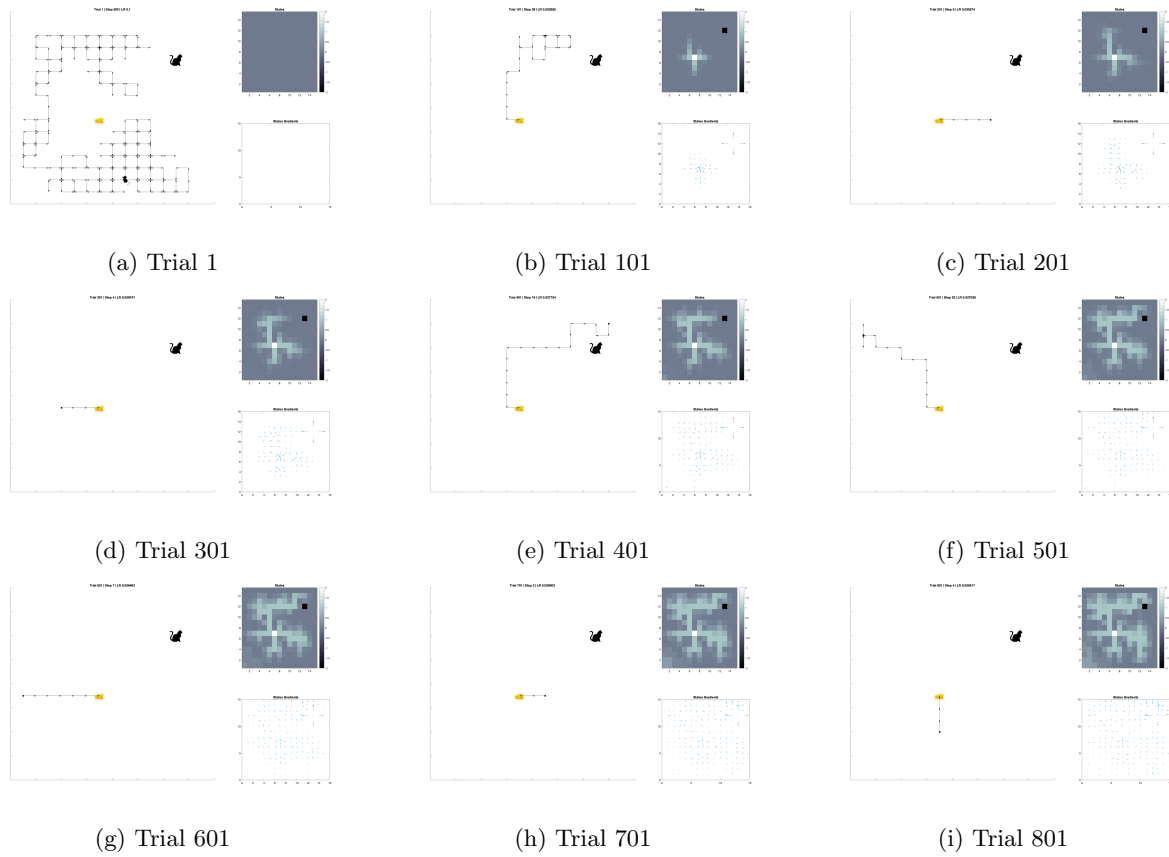


Figure 5: TD Rule - Agent passed way, States values and States gradients in some of the trials. $LR = \frac{0.1}{1+\log t}$, $\gamma = 1$, $\lambda = 0.9$

As can be seen in Figure 5, the agents learn the value of the states over the trials. By comparing Figures 1 and 5 we can infer that TD rule can learn the values of the states faster. This result is more obvious in Figure 6.

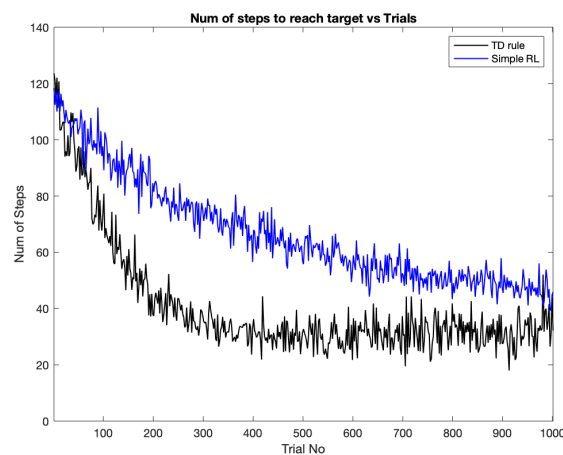


Figure 6: TD-Rule - Number of steps to reach the target vs number of trials of learning. I run the model 200 times and average over all of the vectors of number of the steps in each of the iterations.

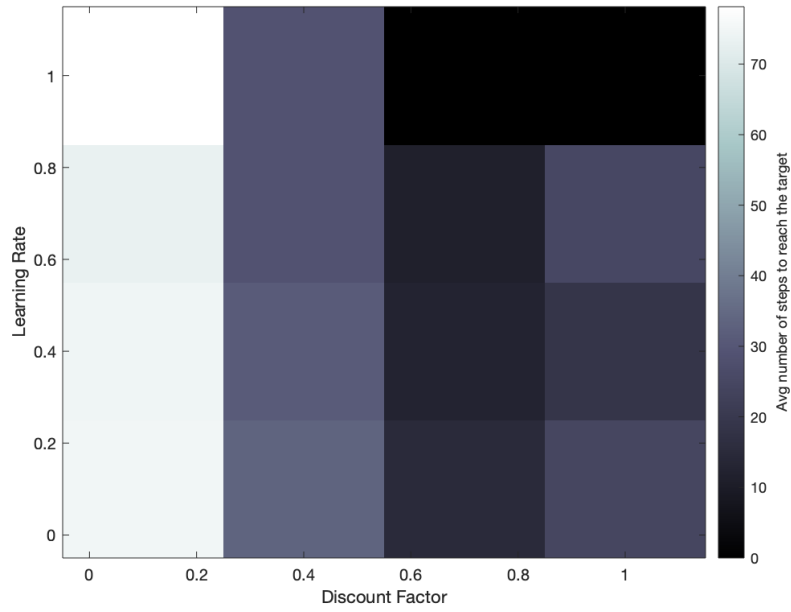


Figure 7: TD-Rule - Average of passed steps to reach one target of the last 50 trials for different learning rates and discounting factors

As can be seen in Figure 7, bigger learning rates or discount factors tend to decrement in the number of passed steps to reach the target. In the following you can see the effect of learning rate and discounting factor on the number of passed steps to reach the target for different RL rules and conditions.

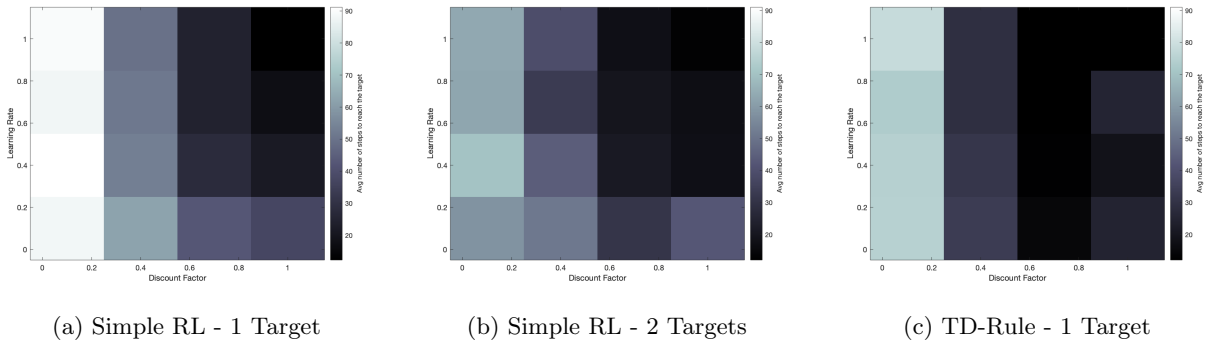


Figure 8: Average of passed steps to reach one target of the last 50 trials for different learning rates and discounting factors during different conditions

It is obvious that TD-rule needs very lower trials to reach the target in comparison with simple RL method and even simple RL with 2 targets! (Figure 8).