

Implementing Actor-Critic Architecture on Grid World

CAP 6629 Reinforcement Learning Project 3

Fanchen Bao

1 Introduction

While tabular method in reinforcement learning can solve problems with a small state and action space, it suffers from the curse of dimensionality when the state or action space becomes large, either due to a lot of discrete choices or being continuous. To address reinforcement learning problems with a large state or action space, one approach is to leverage neural network to model the state or action space. This way, instead of relying on a table to store state information, one can obtain it by passing state as input to a neural network model. Using neural network to model state or action space drastically reduces the burden on memory, and is thus a better method to handle reinforcement learning problems in the real world.

The Actor-Critic architecture is a classic usage of neural network model in reinforcement learning. Generally speaking, the actor is a neural network model of a parameterized policy function, where passing in the current state returns the probability distribution (or density) of all possible actions. The critic is a neural network model of a parameterized state (or action) value function, where passing in the current state (or state-action pair) yields its corresponding state (or action) value. During training, the goal of the actor is to maximize the probability of the action that can maximize the additional return at each state (i.e. advantage, see Section 2 for more detail). The goal of the critic is to align itself with the maximum return at each state. One can view the actor-critic relationship as a competitive one, in which the actor constantly wants to push the return higher than the value estimated by the critic, whereas the critic always wants to catch up with the current return. The two chase each other until the maximum return is approximated by the critic and achieved via the actions predicted by the actor.

The aim of this report is to implement the Actor-Critic architecture for a simple reinforcement learning task. The task chosen is to navigate a grid world, such as the one shown in Figure 1, from top left (start state) to bottom right (end state).

This report is organized in the following manner. Section 2 introduces mathematical notations and concepts of the Actor-Critic architecture. Section 3 provides pseudo-code for the Actor-Critic architecture and describes its implementation details. Section 4 evaluates the performance of the Actor-Critic architecture under different grid worlds. Finally, Section 5 compares the grid world performance under the Actor-Critic architecture with that under the tabular method and concludes the report.

2 Formal Definition

To discuss the mathematical notation of the Actor-Critic architecture, we must start with a parameterized total reward function. In tabular method, there is no total reward function, because it is possible to obtain the discounted total rewards at each step via tabulation. However, this is not possible when the problem involves large or continuous state or action space. Hence, a parameterized total reward function, as defined in Equation (1), is needed.

$$J(\theta) = \sum_{s \in S} d^\pi(s) V^\pi(s) = \sum_{s \in S} d^\pi(s) \sum_{a \in A} \pi_\theta(a|s) Q^\pi(s, a) \quad (1)$$

Here, θ is the trainable parameters. d^π is the stationary distribution of each state under the parameterized policy function π_θ . V^π and Q^π are the state and action value function following π_θ .

To achieve the maximum total reward, we can use gradient ascent technique to tune θ . To use gradient ascent, we need to compute the gradient of $J(\theta)$. According to the policy gradient theorem (see [1] for its mathematical proof), the gradient of $J(\theta)$ can be written as Equation (2).

$$\begin{aligned}\nabla_\theta J(\theta) &= \nabla_\theta \sum_{s \in S} d^\pi(s) \sum_{a \in A} \pi_\theta(a|s) Q^\pi(s, a) \\ &\propto \sum_{s \in S} d^\pi(s) \sum_{a \in A} \nabla_\theta \pi_\theta(a|s) Q^\pi(s, a) \\ &\propto \mathbb{E}_\pi [Q^\pi(s, a) \nabla_\theta \ln \pi_\theta(a|s)]\end{aligned}\tag{2}$$

It is important to note that Equation (2) is not the only way to compute $\nabla_\theta J(\theta)$. According to Schulman [2], $\nabla_\theta J(\theta)$ can also be expressed as Equation (3).

$$\begin{aligned}\nabla_\theta J(\theta) &= \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} (Q^\pi(s_t, a_t) - V^\pi(s_t)) \nabla_\theta \ln \pi_\theta(a|s) \right] \\ &= \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} (G_t - V^\pi(s_t)) \nabla_\theta \ln \pi_\theta(a|s) \right]\end{aligned}\tag{3}$$

Here, we can substitute $Q^\pi(s_t, a_t)$ with G_t , the true discounted reward, because $Q^\pi(s_t, a_t) = \mathbb{E}_\pi[G_t | S_t, A_t]$. The expression $G_t - V^\pi(s_t)$ is called the **advantage**. It describes the difference between the true discounted reward at time t versus the estimated state value.

Equation (4), which is very similar to Equation (3), is the basis for constructing our Actor-Critic architecture.

$$\nabla_\theta J(\theta) = \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} (G_t - V_w(s_t)) \nabla_\theta \ln \pi_\theta(a|s) \right]\tag{4}$$

The only difference between the two is the parameterization of the state value function in Equation (4). We designate $\pi_\theta(a|s)$ as the actor and $V_w(s_t)$ as the critic. They are both approximated by neural network models, parameterized by θ and w , respectively. The reason we choose Equation (4) as the basis for the Actor-Critic architecture is the separation of the actor and critic in terms of training. Since there is no inherent link between the two, they can be trained independently.

However, the antagonistic interaction between the actor and critic is also apparent, because they are linked by the discounted reward G_t . G_t is completely controlled by the actor. The goal of the actor is to predict the best action that can maximize G_t , which in turn also maximizes the advantage. Meanwhile, the goal of the critic is to align itself with the current G_t and reduce the advantage to zero. The antagonistic interaction between the actor and critic can be summarized as follows:

1. The actor increases the advantage.
2. The increased advantage forces the critic to optimize such that the advantage is reduced again.
3. The reduced advantage forces the actor to optimize, and we go back to 1.

This interaction loop does not end until both the actor and critic converge, which leaves us with the parameter θ that maximizes $J(\theta)$. By definition, this will solve any reinforcement learning problem.

3 Implementation

The Actor-Critic architecture is implemented in a grid world as shown in Figure 1.

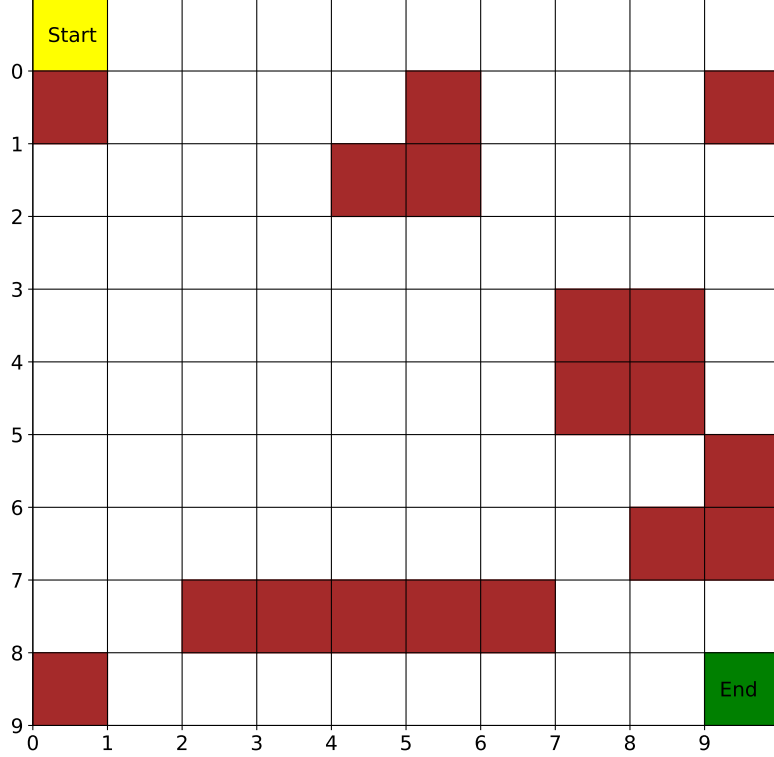


Figure 1: A Grid World with Obstacles

The grid world is 10×10 in size, with the start and end state marked by yellow and green, respectively. The brown cells represent blockages in the grid world where an agent cannot enter. The rule of the grid world is that an agent can only make a move in four directions: up, down, left, and right. Any move incurs a penalty of -1. If an action leads an agent out of bound or into the blockage, the agent remains where it is but still accumulates the penalty. The goal of the grid world is to use the Actor-Critic architecture to learn the optimal path from the start to end state while avoiding all the blockages.

3.1 Pseudo-code

The pseudo-code for the Actor-Critic architecture implementation is provided in Algorithm 1. Intuitively, for each episode, the pseudo-code accumulates C_{val} , P_{act} , and R at each step. Based on R , we can obtain the discounted total reward G at each step. We then use the mean squared error between G and C_{val} as the loss function to tune V_w . This loss function is a reasonable choice because it is widely used when fitting a neural network regressor. That said, it is worth mentioning that the documentation on TensorFlow [3] and Keras [4] use Huber Loss as the critic loss function.

Another important thing to point out is the computation of the actor loss. It is computed as $-\ln(P_{act}) \cdot A$. We need to compute the logarithm of the action probability because Equation

Algorithm 1 Actor-Critic Architecture

```
1: Initialize actor  $\pi_\theta$  with initial random parameters  $\theta$ 
2: Initialize critic  $V_w(s)$  with initial random parameters  $w$ 
3:  $ep \leftarrow 0$  ▷ keep track of the number of episodes experienced
4:  $max\_eps$  ▷ Maximum number of episodes allowed
5:  $T$  ▷ Maximum number of steps per episode
6:  $\alpha_\theta \leftarrow 0.01$  ▷ Learning rate of the actor
7:  $\alpha_w \leftarrow 0.01$  ▷ Learning rate of the critic
8:  $\gamma \leftarrow 0.90$  ▷ Discount rate
9:  $C_{val} \leftarrow []$  ▷ Empty array to store the critic value at each step in an episode
10:  $P_{act} \leftarrow []$  ▷ Empty array to store probability of a sampled action from the actor at each step in an episode
11:  $R \leftarrow []$  ▷ Empty array to store the reward at each step in an episode
12:  $G \leftarrow []$  ▷ Empty array to store the discounted total reward at each step in an episode
13: Initialize starting state  $s$ 
14:
15: repeat
16:   for  $t = 1 \dots T$  do: ▷ Each step in an episode
17:      $a, P_a \leftarrow \pi_\theta(a|s)$  ▷ Sample action and its probability
18:      $c \leftarrow V_w(s)$  ▷ Obtain estimated critic value
19:      $r \leftarrow -1$  ▷ Obtain the reward, which is always -1
20:     Append  $c$  to  $C_{val}$ 
21:     Append  $P_a$  to  $P_{act}$ 
22:     Append  $r$  to  $R$ 
23:     Update  $s$  by applying  $a$  to  $s$ 
24:   end for
25:
26:    $g \leftarrow 0$  ▷ Accumulation of discounted reward
27:   for  $r$  in reversed  $R$  do: ▷ The discounted total reward is computed in reverse
28:      $g \leftarrow r + \gamma g$ 
29:     Append  $g$  to  $G$ 
30:   end for
31:   Reverse  $G$ 
32:
33:    $A \leftarrow G - C_{val}$  ▷ Obtain advantage via pair-wise subtraction
34:    $L_c \leftarrow A^2 / 2T$  ▷ Mean squared error, critic loss
35:    $L_a \leftarrow -\ln(P_{act}) \cdot A$  ▷ Dot product, actor loss
36:
37:   Compute  $\nabla_\theta \pi_\theta$  based on  $L_a$  ▷ Obtain actor gradient
38:    $\theta \leftarrow \theta + \alpha_\theta \nabla_\theta \pi_\theta$  ▷ Update actor parameters
39:   Compute  $\nabla_w V_w$  based on  $L_c$  ▷ Obtain critic gradient
40:    $w \leftarrow w + \alpha_w \nabla_w V_w$  ▷ Update critic parameters
41:
42:    $ep \leftarrow ep + 1$ 
43: until  $ep = max\_eps$ 
44: return  $\pi_\theta$  and  $V_w$  ▷ The trained actor and critic
```

(4) explicitly requires so. We multiply by A , because the advantage gives us the signal whether the actor needs substantial update. If A is small, that means the actor needs to push for a better G_t . Hence it will update itself. If A is large, that means the actor is already producing good actions. Hence it will not make major changes. Finally, we need to add a negative sign, because when the loss function is being minimized during training, we are actually maximizing $\ln(P_{act}) \cdot A$, which is exactly the goal of the actor.

3.2 Neural Network Models

Algorithm 1 is written in Python, relying on Keras for constructing the neural network model for the actor and critic. Diagrams of the neural network models are shown in Figure 2. Both neural networks have two input nodes, corresponding to the row and column index on the grid world. Both also contain a single hidden layer with 128 nodes. The activation function for all hidden nodes is “relu”. The critic has one node in the output layer, with the activation function being “relu”. However, since the reward used in the grid world is negative, a negative sign is added to the critic output. The actor has four nodes in the output layer, corresponding to each action. The activation function for the output layer in the actor is “softmax”. It is worth mentioning that the documentation on TensorFlow [3] and Keras [4] let the actor and critic share the same hidden layer. But in our implementation, we have decided to make their training independent of each other.

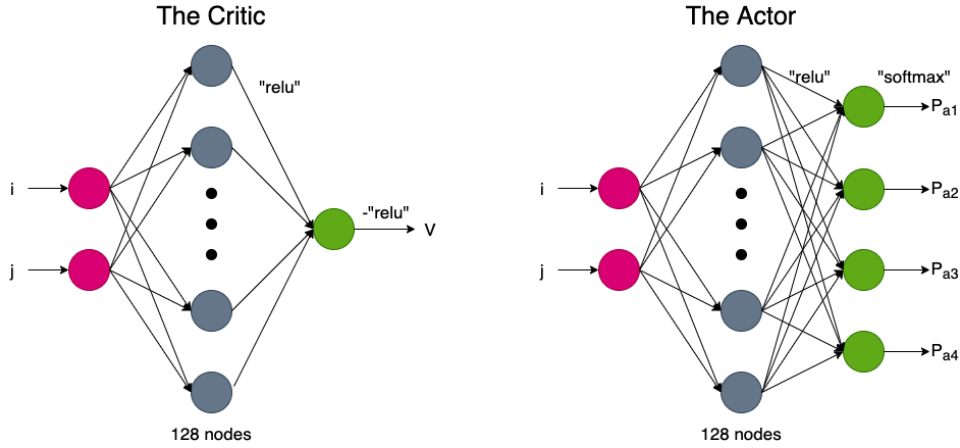


Figure 2: Neural network models for the critic and actor

The gradient of the critic and actor loss function is computed using the `GradientTape` API from TensorFlow. The training consists of 200 episodes, with each episode not exceeding 1000 steps.

4 Evaluation

The evaluation of the Actor-Critic architecture on the grid world is displayed in Figures 3, 4, and 5. These three figures represent three snapshots of the training process at the beginning, middle, and end of the 200 episodes, respectively. In each figure, subplot (A) is the steps-to-go curve at each episode. Since the maximum step per episode is capped at 1000, if an episode fails to reach the end state, the steps-to-go will be shown as 1000. The green horizontal line represents the optimal path with 18 steps. Subplot (B) shows the critic and actor loss in one plot, with the former using the left y-axis and the latter the right. Subplot (C) shows an example path constructed from the currently trained actor. And subplot (D) shows the heat-map of the state values based on the currently trained critic.

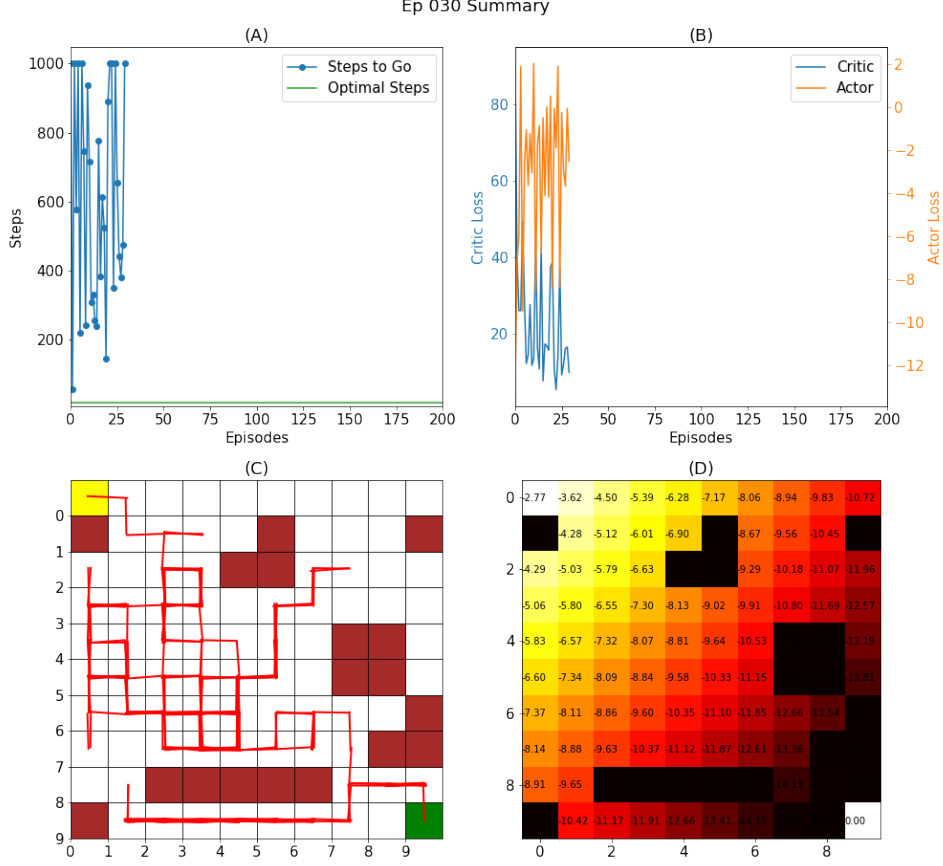


Figure 3: Training summary at 10 episodes, two-opening grid world

There are a few interesting observations. First, quite a few failed episodes occur in the training, even towards the end (e.g. Figure 5). This highlights the stochastic nature of the Actor-Critic architecture, where the action is always sampled from the actor according to their probabilities and there is no guarantee that the action with the highest probability is always going to be selected. That said, if we do not pay attention to the failed episodes, the steps-to-go seems to converge towards the optimal path (Figure 5).

Second, there is a lot of variation in both loss functions (subplot (B)), yet the critic loss seems to gradually decrease and the actor loss seems to converge towards 0 (Figure 5). This is understandable. For the critic loss, getting smaller means the critic is better at estimating the true discounted reward. For the action loss, converging towards 0 means the probability of the actor's predicted action is approaching 1, which means the actor is more confident in its choice.

Third, the agent apparently explores much more at the early stage of training. This can be seen in subplot (C) of Figure 3 where the agent explores many parts of the grid world. This is within expectation as the actor has not yet been solidified. This allows the agent the freedom to explore the grid world via random walk. Random walk is crucial to training, because the rule of the grid world mandates that the only way for an agent to reach the end state initially is via random walk (i.e. there is no other incentive or cues to guide it there). If a different reward system is given, where intermediate goals are given to guide the agent through the highly indeterministic area of the grid world, there will surely be much less random walk.

Finally, an obvious shift in the state value can be observed in subplot (D) across the three figures. In Figure 3, the high state value region (light color) is near the starting point; it rotates to the bottom left in Figure 4; and finally lands on the end state in Figure 5. This is a good representation of how the Actor-Critic architecture works out the problem. Initially, especially

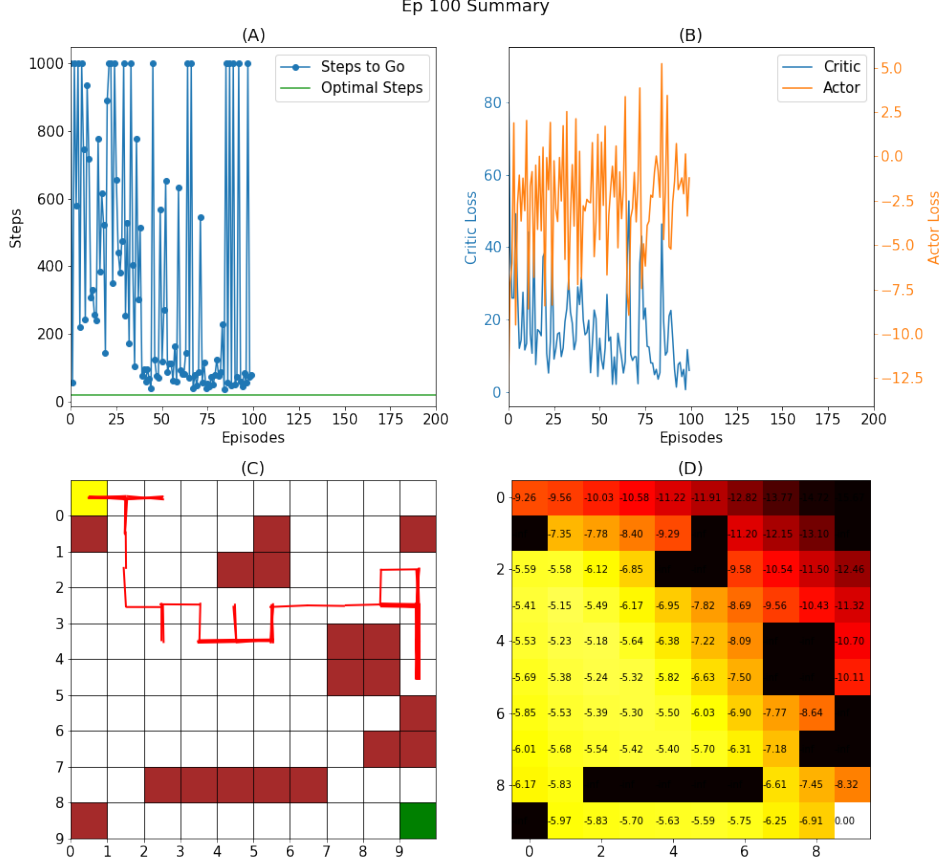


Figure 4: Training summary at 100 episodes, two-opening grid world

before the end state can be reached via random walk, the critic evaluates each state mainly based on how accessible it is. If a state is accessible, i.e. it can be reached from many sides, leaving that state for another will incur fewer penalties, because there is little chance the agent will get stuck. This is the case with the states near the upper left corner. On the contrary, if a state is surrounded by blockages, moving away from that state will likely get the agent stuck and incur a lot of penalties. Hence the states in the bottom right corner have low state values (dark color) at the beginning.

However, the situation starts to change once the end state is reached, because the states leading to the end state will now accumulate the least penalties. From Figure 4, it is apparent that the state values towards the bottom right corner are significantly higher than before (light color). Meanwhile, the low state value region (dark color) moves to the upper right corner as it is a region with good amount of blockages and also quite far from the end state.

Finally, when the critic stabilizes in Figure 5, the high state value region shifts to the bottom right corner. We can even pick out a path based on the color gradient in the heat-map. The top right corner still has the lowest state values (dark color) because it is far from the end state and surrounded by blockages. The top left corner is also low in state value (dark color) because it is also far from the end state. Yet, since there are fewer blockages in the top left corner, its state value is generally higher than the top right corner. The bottom left corner has similar level of state value as the top left. Despite the closeness of the bottom left corner to the end state, the many blockages around it counter this benefit.

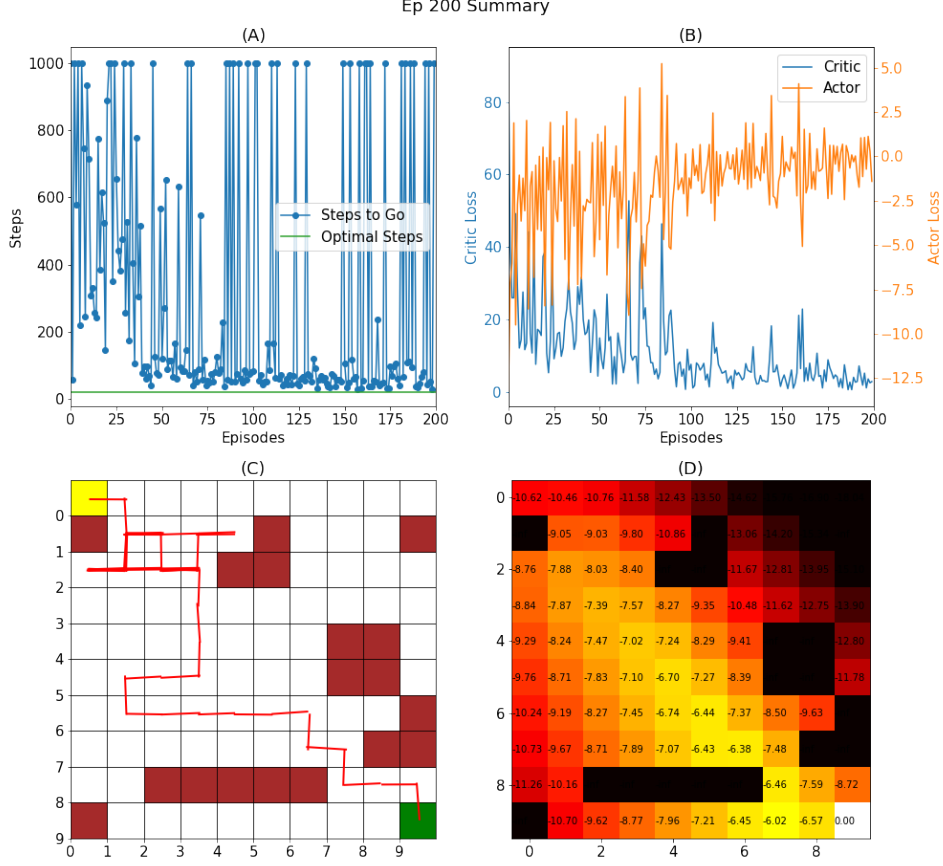


Figure 5: Training summary at 200 episodes, two-opening grid world

4.1 The One-opening Grid World

The evaluation shown above is based on a two-opening grid world. However, it is interesting to report that among all the paths the agent takes during the training process, it almost exclusively uses the right opening. One reason is that, compared to the left opening, the right opening is closer to the center of the grid world, which means it is more likely for the agent to hit it via random walk, especially during the early stage of training. Another reason is that the right opening is much closer to the end state than the left opening, which naturally gives it more importance, especially towards the end of training.

This observation begs the question: will an agent using the same Actor-Critic architecture be able to reach the end state if the right opening is blocked? To answer this question, we create a second grid world that has only one opening and run the same algorithm on it. For brevity, only the summary of the final episode is shown in Figure 6.

As demonstrated in subplot (C) and (D), the algorithm is not able to find a path through the left opening. The reason for such failure can be sourced from subplot (A) where we can clearly see that there are very few cases where the agent reaches the end state (i.e. successful episodes). Since only the successful episodes provide training information regarding the whereabouts of the end state and left opening, lack of successful episodes means that the actor and critic are under-trained. It is thus not surprising that the algorithm fails to find the correct path.

A deeper question to ask is why only few successful episodes are achieved in Figure 6. The answer is the same as why the right opening is preferred: it is simply not likely for an agent to reach the left opening via random walk. The more often the agent fails to reach the left opening during the early stage of training, the more likely it will not get there in the later stage as the actor and critic are solidified by the training data from the other states. Therefore, it is not

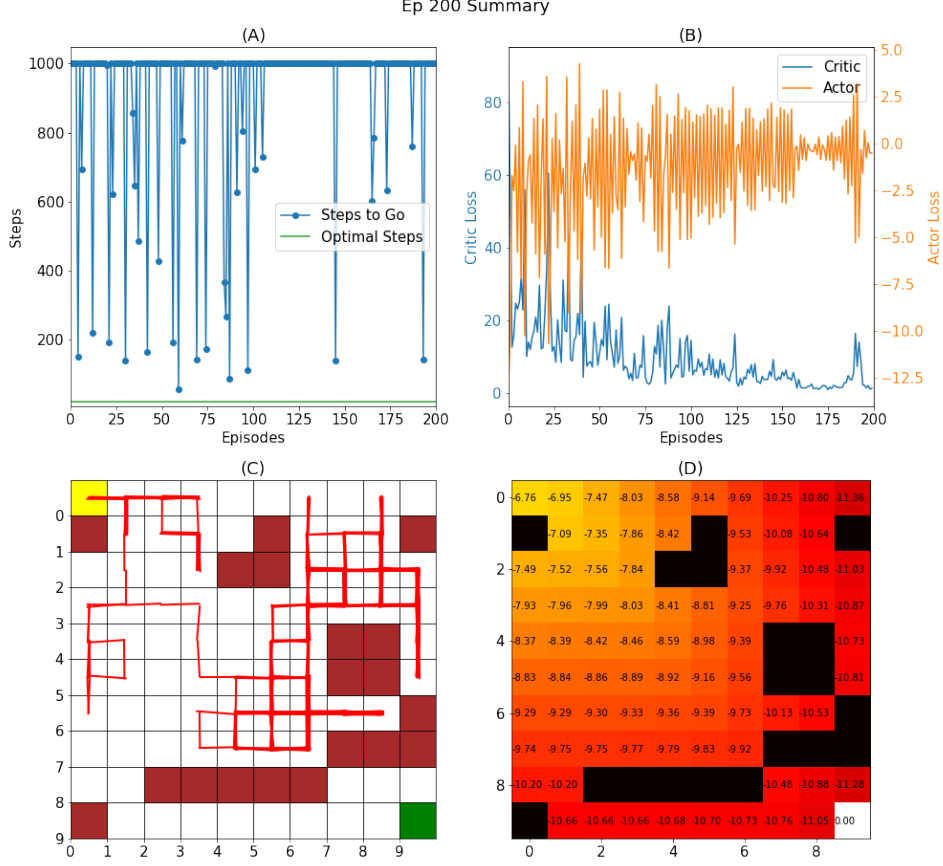


Figure 6: Training summary at 200 episodes, one-opening grid world, top left to bottom right

a stretch to say that given the current algorithm, we cannot solve the one-opening grid world problem.

4.2 Or can we?

The main question is how to get the agent to the left opening via random walk. It is apparently not likely when the agent starts from the top left corner, since there are too many distractions. But how about starting from the bottom right corner? There is no rule preventing us from using the bottom right corner as the starting point. A path found from bottom right to top left is equivalent from a path found from top left to bottom right. Therefore, it is perfectly fine going from the end state backwards toward the start state.

Yet, the benefit of doing so is tremendous. Since the bottom right and bottom left corners are heavily blocked, there is little ambiguity of where to go. With such reduced distraction, it is expected that an agent will be able to reach the left opening easily. Once the left opening is reached, finding the start state shouldn't be too much trouble, since it is located in an accessible area.

To confirm our hypothesis, we run the algorithm on the one-opening grid world backwards, with the start and end state swapped. The results are shown in Figures 7, 8, and 9, representing the beginning, middle, and end of the training process, respectively.

Compared to the one-opening grid world going from top left to bottom right, the reversed direction yields two major benefits. To begin with, at all three training stages, the agent is able to find its way to the start state (see subplot (C) in all three figures). This confirms our hypothesis that with a highly deterministic path from the end state to the left opening, the

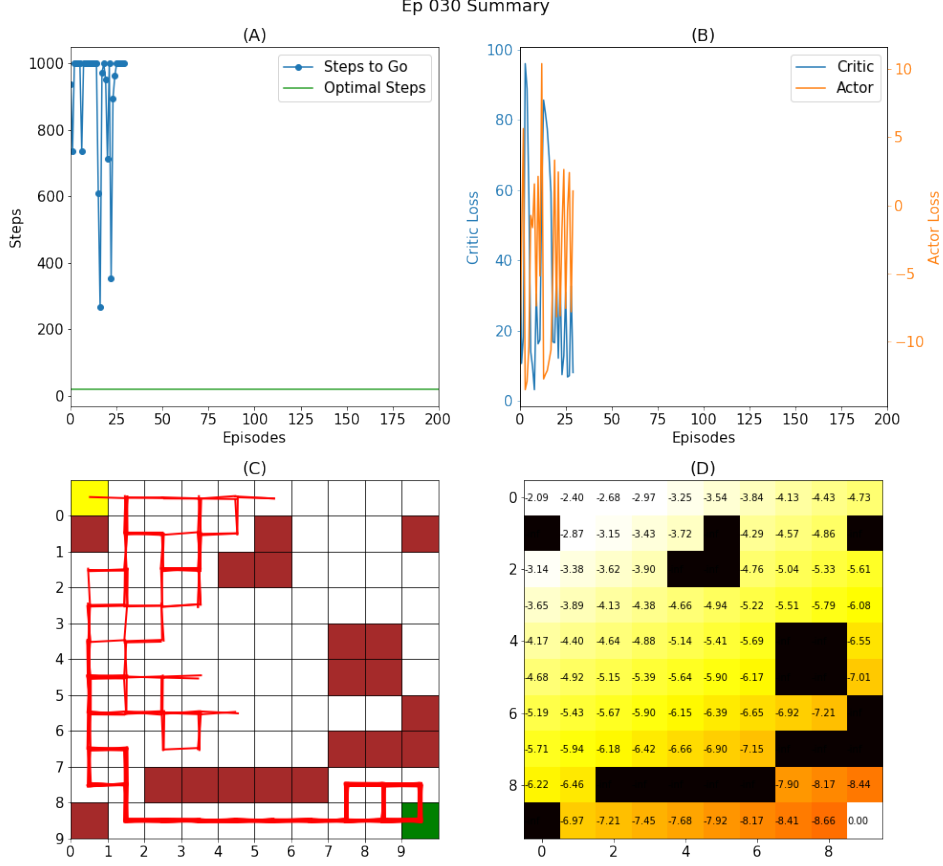


Figure 7: Training summary at 30 episodes, one-opening grid world, bottom right to top left

agent is able to first locate the left opening, and then random walk its way to the start state.

Second, from subplot (A), it is clear that going from the end to the start state results in many more successful episodes than the other way around. This provides the actor and critic more training samples, which further increases the agent’s chance of having a successful episode. This virtuous cycle is in stark contrast to the vicious cycle when the agent is trained to go from the start to the end state.

Although the path obtained in 200 episodes is far from optimal, swapping the start and end state definitely solves a seemingly unsolvable reinforcement learning problem, without any modification to the core algorithm of the Actor-Critic architecture. It is highly likely that had we extended the training episodes, the reversed method would be able to converge its steps-to-go curve closer to the optimal value.

A good lesson we can learn from this example is that domain knowledge in machine learning is very important. With regard to the grid world, our domain knowledge points out that the only opening in the grid world is closer to the end state and that the path from that opening to the end state is deterministic. This domain knowledge allows us to swap the start and end state, and make the problem easier to solve. There are, of course, many other ways to break down the problem (e.g. set up intermediate goal at the left opening), but the key point is to always remember that one small piece of domain knowledge can be worth hundreds of hours of CPU/GPU time in machine learning.

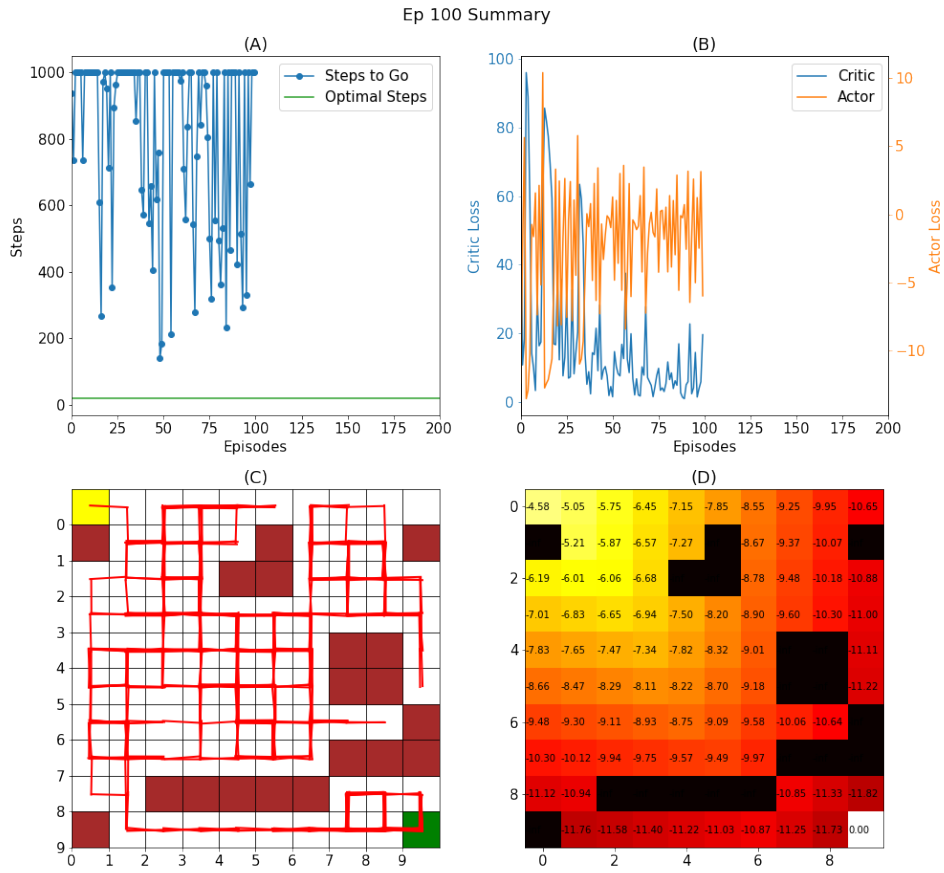


Figure 8: Training summary at 100 episodes, one-opening grid world, bottom right to top left

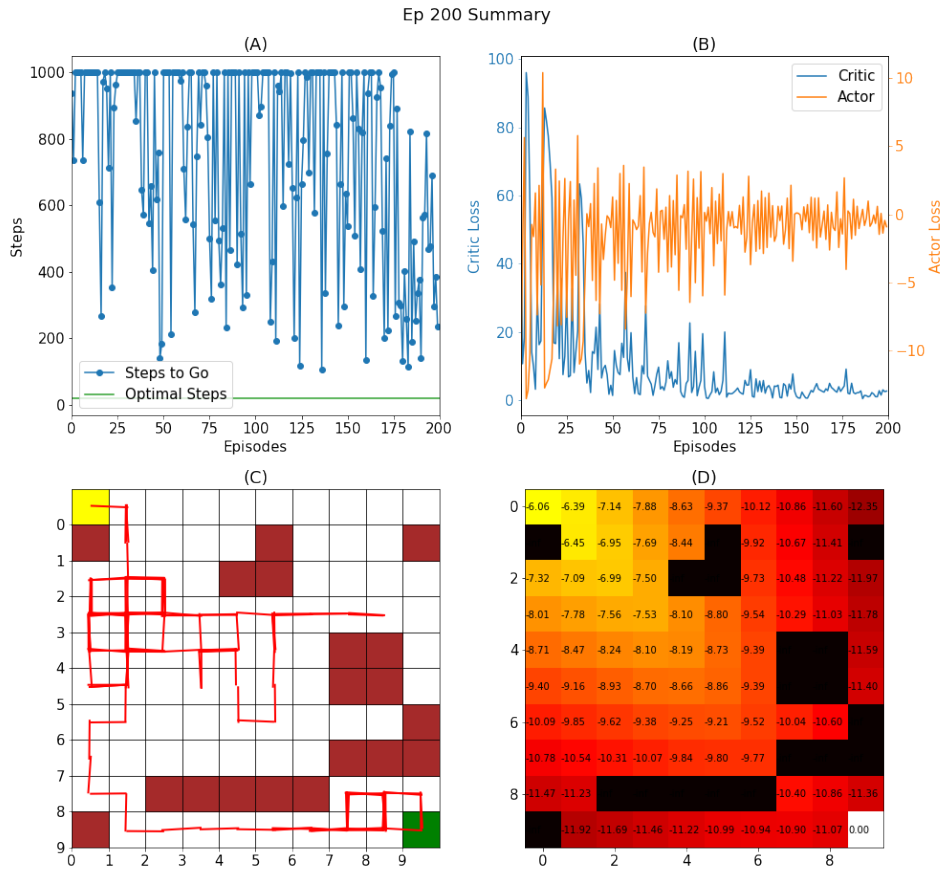


Figure 9: Training summary at 200 episodes, one-opening grid world, bottom right to top left

4.3 Comparison with Project 2

A quick comparison between the performance of the Actor-Critic architecture and that of the TD- λ tabular method from Project 2 (Figure 10) shows that the tabular method is much more stable in the training process (i.e. less fluctuation in steps-to-go), produces deterministic path, and normally converges its solution in shorter period of time. This tells us that the Actor-Critic architecture is not the right tool for use cases where the state and action space is small. It will be interesting to investigate at what magnitude of state and action space is the Actor-Critic architecture a better option than the tabular method.

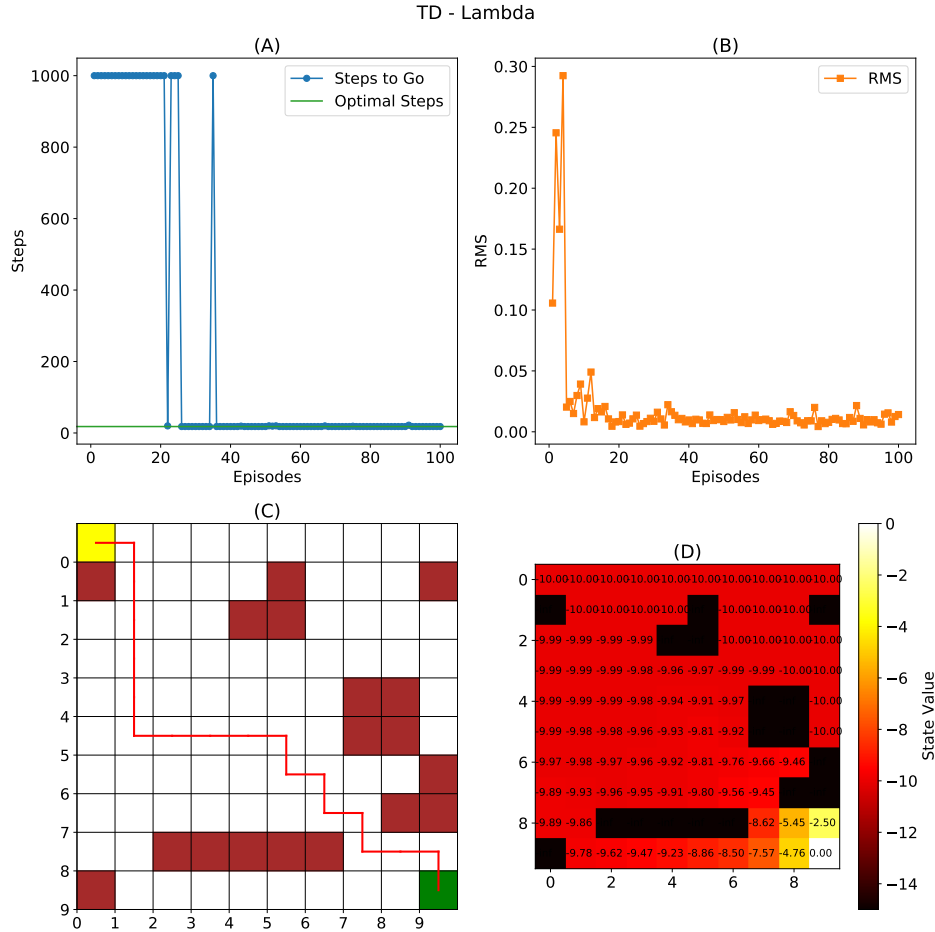


Figure 10: Grid world solved by TD- λ tabular method

5 Discussion And Conclusion

In this report, we have explained the mathematical reasoning behind the Actor-Critic architecture, showed its pseudo-code, and observed the performance of the Actor-Critic architecture over two types of grid world. We have found that our implementation of the actor and critic depends heavily on random walk to reach the end state at the beginning stage of training. This means if the end state or any bottleneck state (e.g. the left opening) is located in a region not easily accessible by random walk, it will be very difficult for the agent to obtain sufficient knowledge about them and the tuning of the actor and critic will fail.

One way to resolve this issue is to leverage the domain knowledge. In the case of the one-opening grid world, swapping the start and end state has significantly simplified the problem and made it solvable.

Finally, we compare the performance of solving the grid world problem using the Actor-Critic architecture and the TD- λ tabular method. The result shows that the tabular method has an advantage over the Actor-Critic architecture when the state and action space in the reinforcement learning problem is small.

References

- [1] L. Weng, “Policy Gradient Algorithms,” Apr. 2018. [Online]. Available: <https://lilianweng.github.io/2018/04/08/policy-gradient-algorithms.html>
- [2] J. Schulman, P. Moritz, S. Levine, M. Jordan, and P. Abbeel, “High-Dimensional Continuous Control Using Generalized Advantage Estimation,” *arXiv:1506.02438 [cs]*, Oct. 2018, arXiv: 1506.02438. [Online]. Available: <http://arxiv.org/abs/1506.02438>
- [3] “Playing CartPole with the Actor-Critic Method | TensorFlow Core,” Apr. 2021. [Online]. Available: https://www.tensorflow.org/tutorials/reinforcement_learning/actor_critic
- [4] A. Nandan, “Keras documentation: Actor Critic Method,” May 2020. [Online]. Available: https://keras.io/examples/rl/actor_critic_cartpole/