



# Implementing Actor-Critic Architecture on Grid World

CAP6629 Reinforcement Learning Project 3

Fanchen Bao



# Introduction

- Tabular method
  - Small state and action space
  - Curse of dimensionality
- Actor-Critic
  - Neural network to simulate state value function and policy function
  - Good for large or continuous state or action space



## Formal Definition

- Parameterize total reward function (and policy function)

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

## Formal Definition

- Policy gradient theorem [1]

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \nabla_{\theta} \sum_{s \in \mathcal{S}} d^{\pi}(s) \sum_{a \in \mathcal{A}} \pi_{\theta}(a|s) Q^{\pi}(s, a) \\ &\propto \sum_{s \in \mathcal{S}} d^{\pi}(s) \sum_{a \in \mathcal{A}} \nabla_{\theta} \pi_{\theta}(a|s) Q^{\pi}(s, a) \\ &\propto \mathbb{E}_{\pi} \left[ \underbrace{Q^{\pi}(s, a)}_{\text{Critic}} \nabla_{\theta} \ln \underbrace{\pi_{\theta}(a|s)}_{\text{Actor}} \right]\end{aligned}$$

## Formal Definition

- Policy gradient theorem [2]

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

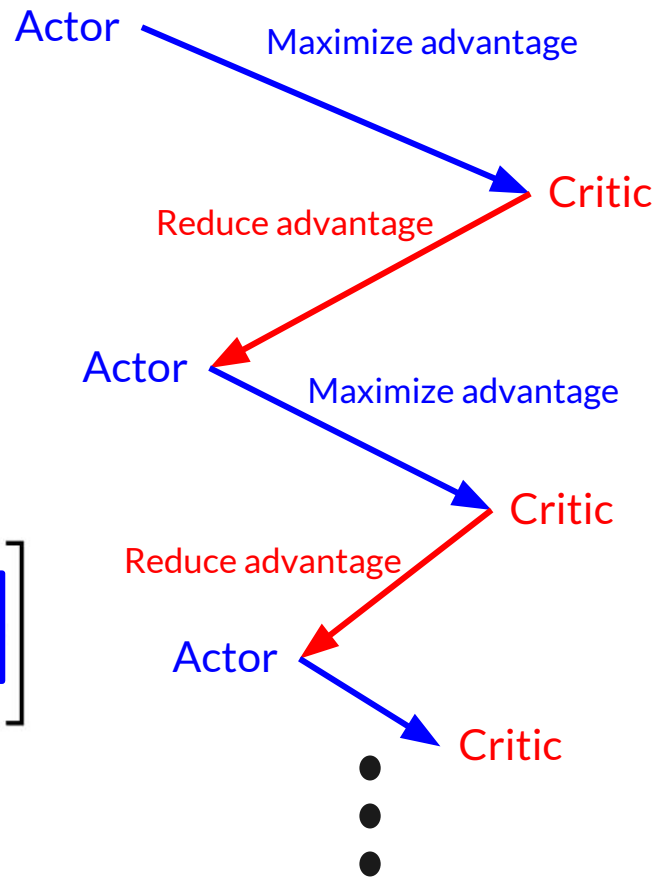
The diagram illustrates the Policy Gradient Theorem. The equation is  $\nabla_{\theta} J(\theta) = \mathbb{E}_{\pi} \left[ \sum_{t=0}^{\infty} (G_t - V_w(s_t)) \nabla_{\theta} \ln \pi_{\theta}(a|s) \right]$ . The term  $(G_t - V_w(s_t))$  is enclosed in a green box and labeled "Advantage" in green text above it. Within this green box, the term  $V_w(s_t)$  is further enclosed in a red box and labeled "Critic" in red text below it. The term  $\pi_{\theta}(a|s)$  is enclosed in a blue box and labeled "Actor" in blue text below it.

## Formal Definition

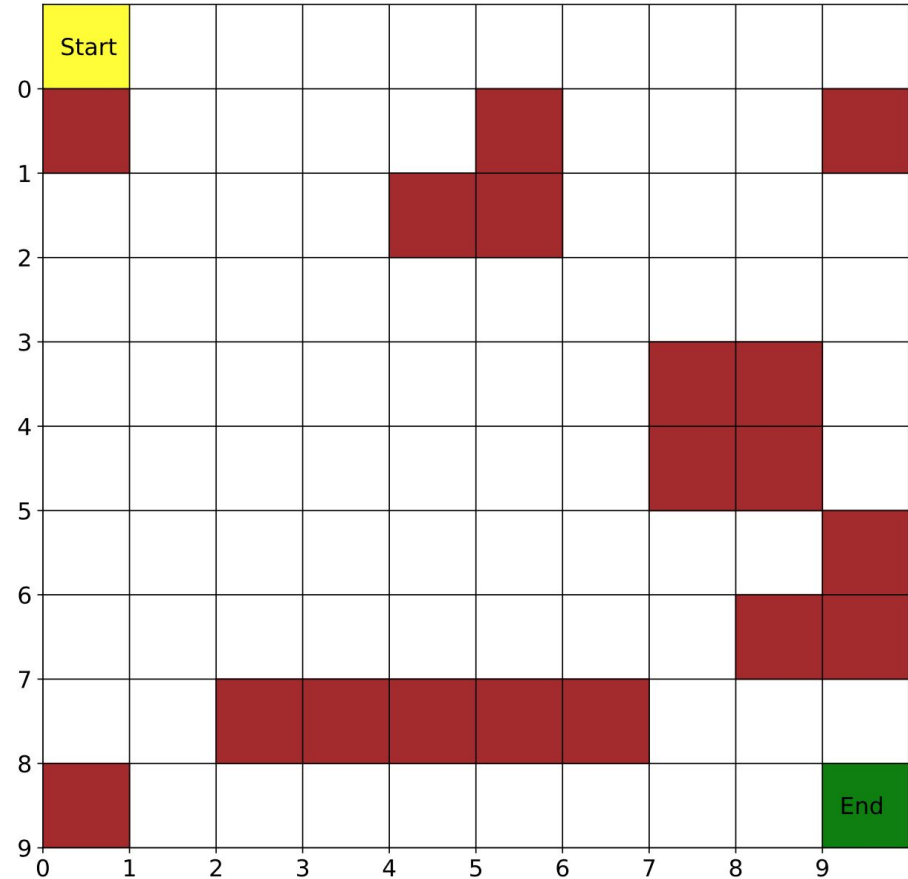
- Antagonistic interaction between actor and critic

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

Critic Actor



- 



# Implementation

- Pseudo-code

- $\gamma = 0.9$
- $\alpha_{\theta} = 0.01$
- $\alpha_w = 0.01$

```
15: repeat
16:   for  $t = 1 \dots T$  do:
17:      $a, P_a \leftarrow \pi_{\theta}(a|s)$ 
18:      $c \leftarrow V_w(s)$ 
19:      $r \leftarrow -1$ 
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:    $g \leftarrow 0$ 
27:   for  $r$  in reversed  $R$  do:
28:      $g \leftarrow r + \gamma g$ 
29:     Append  $g$  to  $G$ 
30:   end for
31:   Reverse  $G$ 
32:
33:    $A \leftarrow G - C_{val}$ 
34:    $L_c \leftarrow A^2 / 2T$ 
35:    $L_a \leftarrow -\ln(P_{act}) \cdot A$ 
36:
37:   Compute  $\nabla_{\theta} \pi_{\theta}$  based on  $L_a$ 
38:    $\theta \leftarrow \theta + \alpha_{\theta} \nabla_{\theta} \pi_{\theta}$ 
39:   Compute  $\nabla_w V_w$  based on  $L_c$ 
40:    $w \leftarrow w + \alpha_w \nabla_w V_w$ 
41:
42:    $ep \leftarrow ep + 1$ 
43: until  $ep = \text{max\_eps}$ 
```

## Monte Carlo

- ▷ Each step in an episode
- ▷ Sample action and its probability
- ▷ Obtain estimated critic value
- ▷ Obtain the reward, which is always -1

- ▷ Accumulation of discounted reward
- ▷ The discounted total reward is computed in reverse

- ▷ Obtain advantage via pair-wise subtraction
- ▷ Mean squared error, critic loss
- ▷ Dot product, actor loss

- ▷ Obtain actor gradient
- ▷ Update actor parameters
- ▷ Obtain critic gradient
- ▷ Update critic parameters



# Implementation

- Pseudo-code

- $\gamma = 0.9$
- $\alpha_\theta = 0.01$
- $\alpha_w = 0.01$

```
15: repeat
16:   for  $t = 1 \dots T$  do:
17:      $a, P_a \leftarrow \pi_\theta(a|s)$ 
18:      $c \leftarrow V_w(s)$ 
19:      $r \leftarrow -1$ 
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:    $g \leftarrow 0$ 
27:   for  $r$  in reversed  $R$  do:
28:      $g \leftarrow r + \gamma g$ 
29:     Append  $g$  to  $G$ 
30:   end for
31:   Reverse  $G$ 
32:
33:    $A \leftarrow G - C_{val}$ 
34:    $L_c \leftarrow A^2 / 2T$ 
35:    $L_a \leftarrow -\ln(P_{act}) \cdot A$ 
36:
37:   Compute  $\nabla_\theta \pi_\theta$  based on  $L_a$ 
38:    $\theta \leftarrow \theta + \alpha_\theta \nabla_\theta \pi_\theta$ 
39:   Compute  $\nabla_w V_w$  based on  $L_c$ 
40:    $w \leftarrow w + \alpha_w \nabla_w V_w$ 
41:
42:    $ep \leftarrow ep + 1$ 
43: until  $ep = \text{max\_eps}$ 
```

- ▷ Each step in an episode
- ▷ Sample action and its probability
- ▷ Obtain estimated critic value
- ▷ Obtain the reward, which is always -1

- ▷ Accumulation of discounted reward
- ▷ The discounted total reward is computed in reverse

## Discounted rewards

- ▷ Obtain advantage via pair-wise subtraction
- ▷ Mean squared error, critic loss
- ▷ Dot product, actor loss
- ▷ Obtain actor gradient
- ▷ Update actor parameters
- ▷ Obtain critic gradient
- ▷ Update critic parameters

# Implementation

- Pseudo-code

- $\gamma = 0.9$
- $\alpha_\theta = 0.01$
- $\alpha_w = 0.01$

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

```

15: repeat
16:   for  $t = 1 \dots T$  do:
17:      $a, P_a \leftarrow \pi_\theta(a|s)$ 
18:      $c \leftarrow V_w(s)$ 
19:      $r \leftarrow -1$ 
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:    $g \leftarrow 0$ 
27:   for  $r$  in reversed  $R$  do:
28:      $g \leftarrow r + \gamma g$ 
29:     Append  $g$  to  $G$ 
30:   end for
31:   Reverse  $G$ 
32:
33:    $A \leftarrow G - C_{val}$ 
34:    $L_c \leftarrow A^2 / 2T$ 
35:    $L_a \leftarrow -\ln(P_{act}) \cdot A$ 
36:
37:   Compute  $\nabla_\theta \pi_\theta$  based on  $L_a$ 
38:    $\theta \leftarrow \theta + \alpha_\theta \nabla_\theta \pi_\theta$ 
39:   Compute  $\nabla_w V_w$  based on  $L_c$ 
40:    $w \leftarrow w + \alpha_w \nabla_w V_w$ 
41:
42:    $ep \leftarrow ep + 1$ 
43: until  $ep = \text{max\_eps}$ 

```

- ▷ Each step in an episode
- ▷ Sample action and its probability
- ▷ Obtain estimated critic value
- ▷ Obtain the reward, which is always -1

- ▷ Accumulation of discounted reward
- ▷ The discounted total reward is computed in reverse

- ▷ Obtain advantage via pair-wise subtraction
- ▷ Mean squared error, critic loss
- ▷ Dot product, actor loss

Compute loss

- ▷ Obtain actor gradient
- ▷ Update actor parameters
- ▷ Obtain critic gradient
- ▷ Update critic parameters

# Implementation

- Pseudo-code

- $\gamma = 0.9$
- $\alpha_{\theta} = 0.01$
- $\alpha_w = 0.01$

```
15: repeat
16:   for  $t = 1 \dots T$  do:
17:      $a, P_a \leftarrow \pi_{\theta}(a|s)$ 
18:      $c \leftarrow V_w(s)$ 
19:      $r \leftarrow -1$ 
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:    $g \leftarrow 0$ 
27:   for  $r$  in reversed  $R$  do:
28:      $g \leftarrow r + \gamma g$ 
29:     Append  $g$  to  $G$ 
30:   end for
31:   Reverse  $G$ 
32:
33:    $A \leftarrow G - C_{val}$ 
34:    $L_c \leftarrow A^2 / 2T$ 
35:    $L_a \leftarrow -\ln(P_{act}) \cdot A$ 
36:
37:   Compute  $\nabla_{\theta} \pi_{\theta}$  based on  $L_a$ 
38:    $\theta \leftarrow \theta + \alpha_{\theta} \nabla_{\theta} \pi_{\theta}$ 
39:   Compute  $\nabla_w V_w$  based on  $L_c$ 
40:    $w \leftarrow w + \alpha_w \nabla_w V_w$ 
41:
42:   ep  $\leftarrow$  ep + 1
43: until ep = max_eps
```

- ▷ Each step in an episode
- ▷ Sample action and its probability
- ▷ Obtain estimated critic value
- ▷ Obtain the reward, which is always -1

- ▷ Accumulation of discounted reward
- ▷ The discounted total reward is computed in reverse

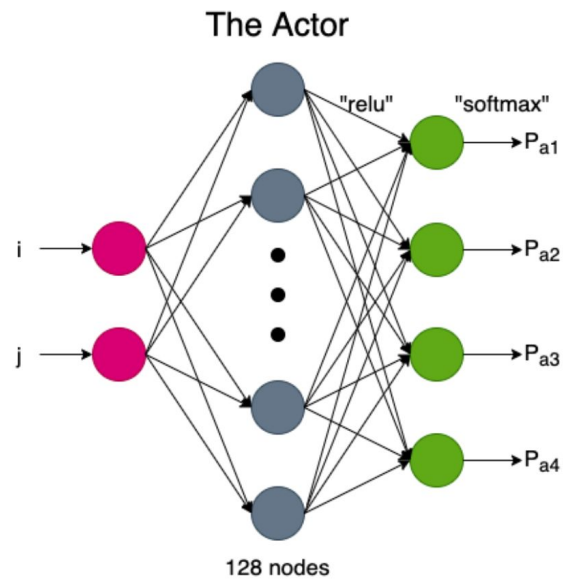
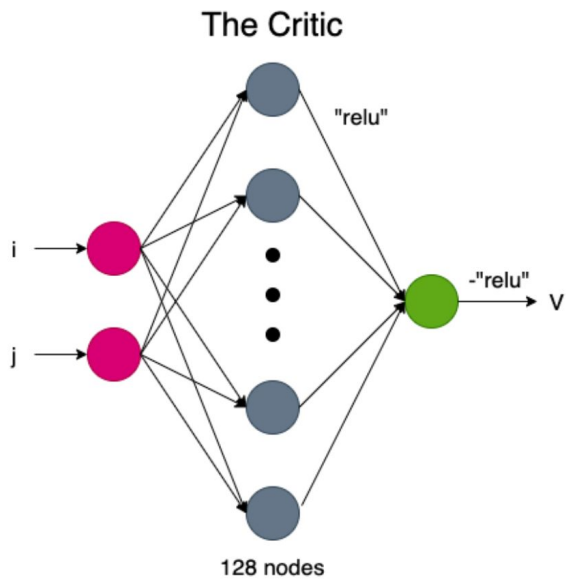
- ▷ Obtain advantage via pair-wise subtraction
- ▷ Mean squared error, critic loss
- ▷ Dot product, actor loss

- ▷ Obtain actor gradient
- ▷ Update actor parameters
- ▷ Obtain critic gradient
- ▷ Update critic parameters

Gradient descent

# Implementation

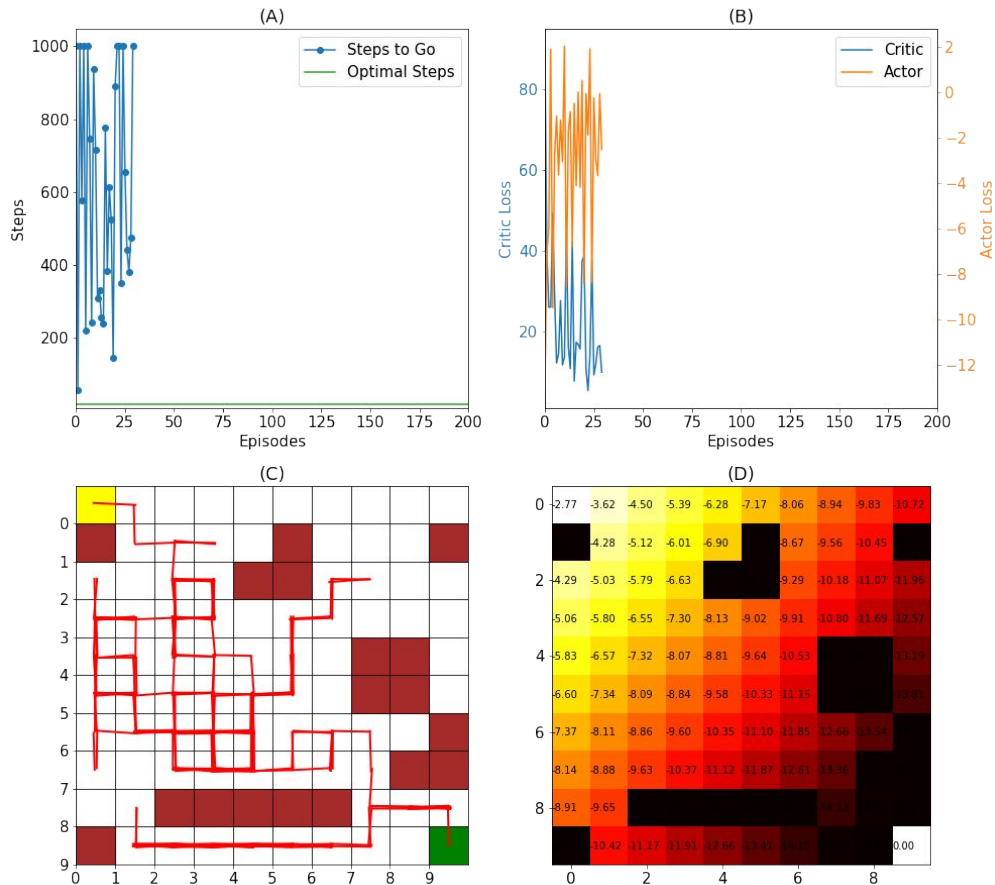
- Neural network



# Evaluation

- Two openings

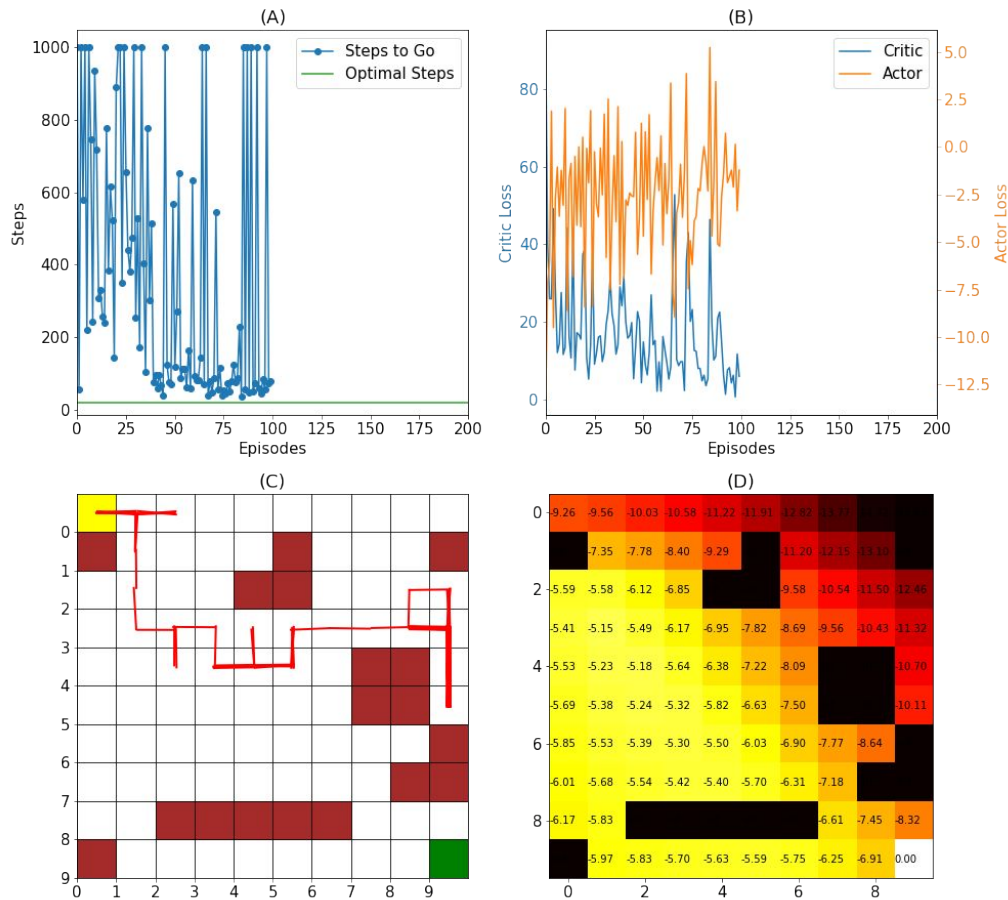
Ep 030 Summary



# Evaluation

- Two openings

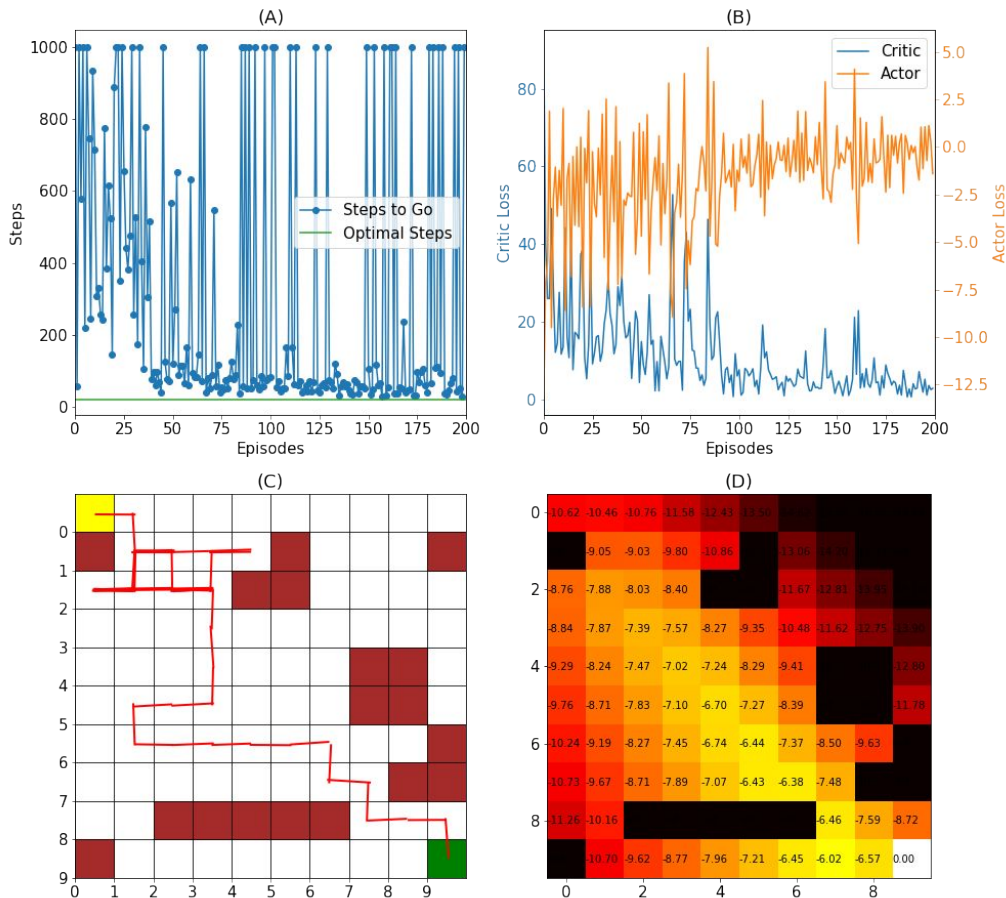
Ep 100 Summary



# Evaluation

- Two openings

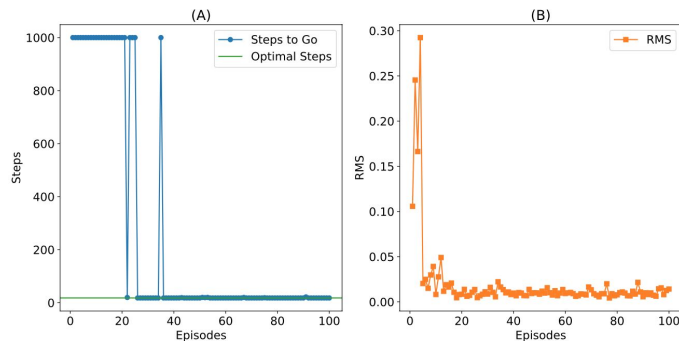
Ep 200 Summary



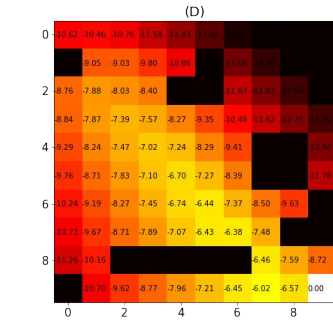
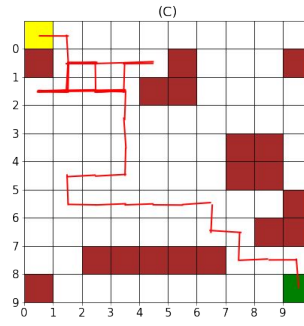
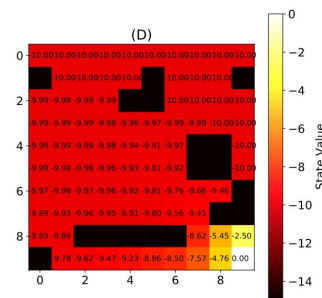
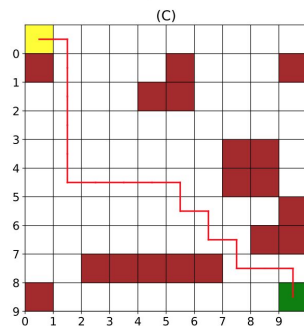
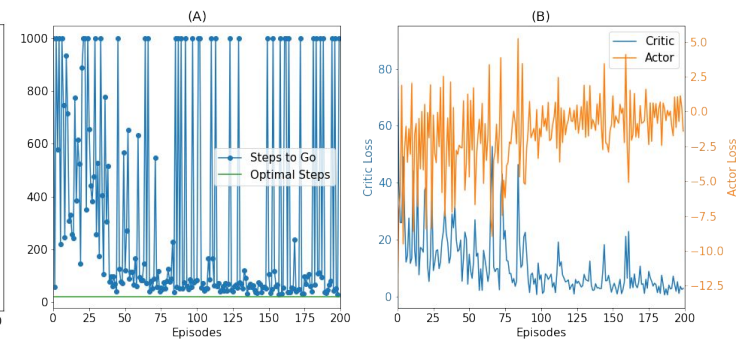
# Evaluation

- Comparison with TD- $\lambda$

TD- $\lambda$



Actor-Critic

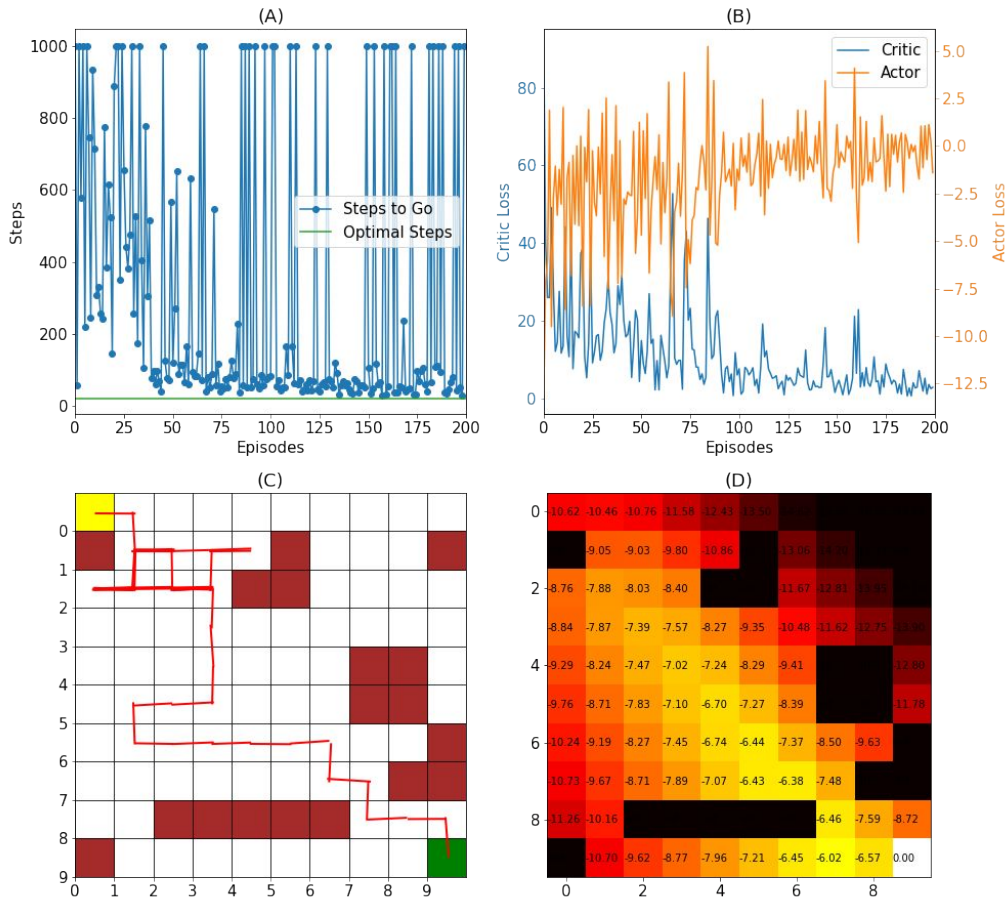




# Evaluation

- Only the right opening is used

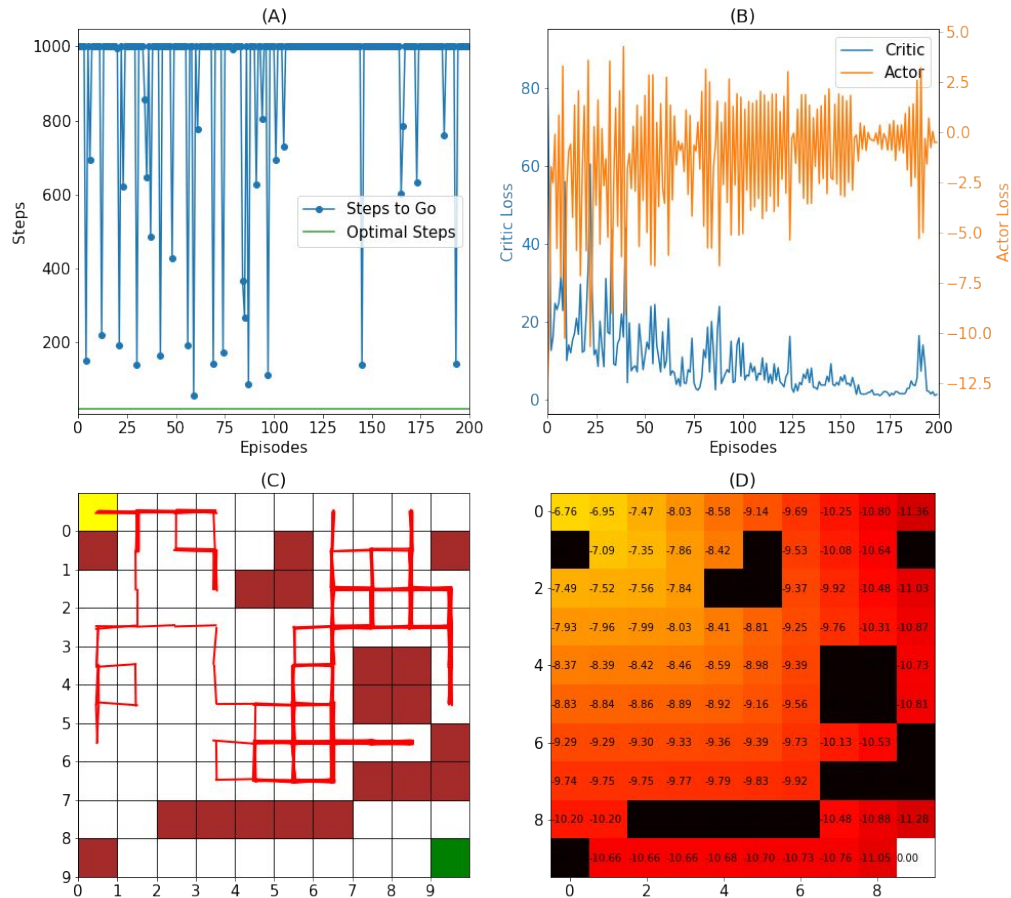
Ep 200 Summary



# Evaluation

- One opening

Ep 200 Summary

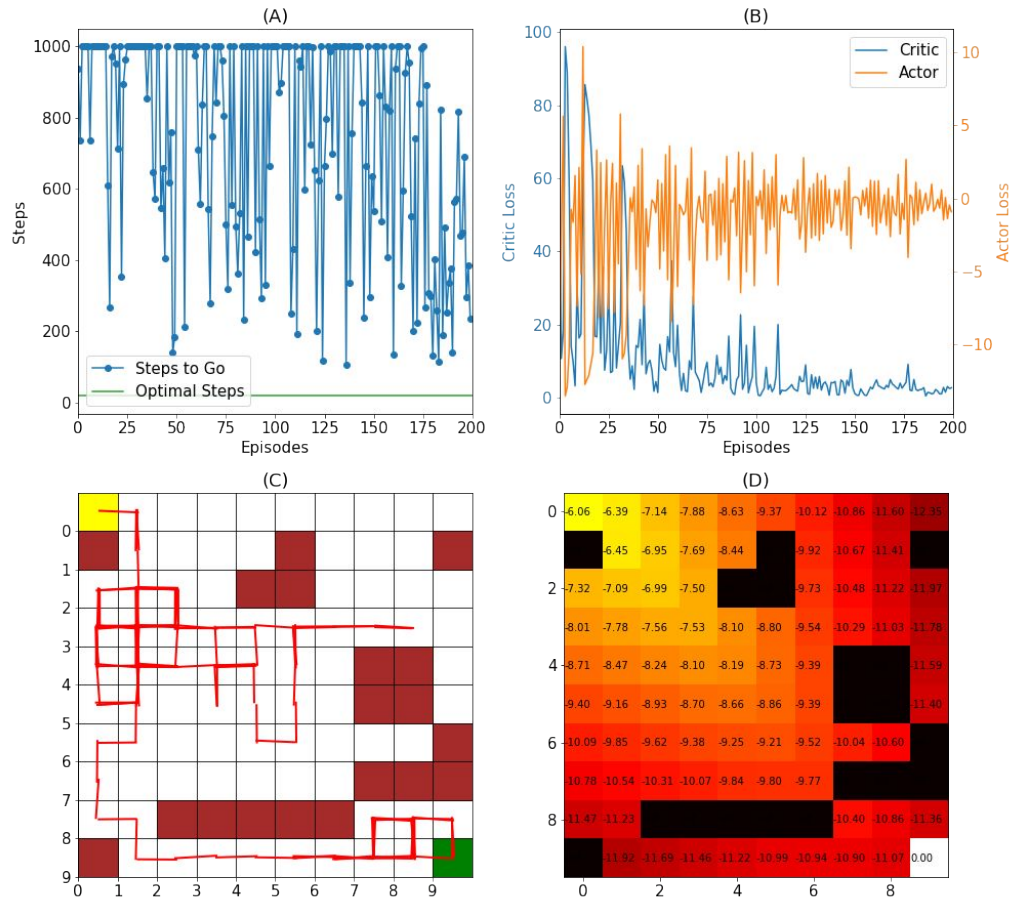


# Evaluation

- One opening
- Bottom right to top left

**Domain knowledge is important!**

Ep 200 Summary





## Conclusions

- Actor-Critic depends on random walk initially.
- Actor-Critic only approximate optimal solution.
- Actor-Critic is slower and less stable in training than TD- $\lambda$ .
- Domain knowledge crucial in solving reinforcement learning problem.



## 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>