

# Tabular Methods - Policy Optimization with On-Policy First-Visit Monte Carlo Control, On-Policy SARSA, and Off-Policy Q-Learning

Ishaan Kandamuri

## I. Introduction

For this assignment, three different model-free RL algorithms: on-policy first visit MC control, SARSA, and Q-learning were implemented to optimize a policy  $\pi$  for the Gymnasium FrozenLake-v1 environment. The task centered around navigating a 4x4 grid, where each tile can be either solid ice, a hazardous hole, the starting point, or the goal. The agent always begins at the top-left corner (0,0) and aims to reach the bottom-right corner (3,3) by moving either left, right, up, or down across 16 possible tiles. This environment is modeled as a Markov Decision Process (MDP). When `is_slippery` is set to true, the intended direction of motion occurs with probability  $\frac{1}{3}$ , while the other  $\frac{2}{3}$  of the time the agent slides into one of the perpendicular directions, each with probability  $\frac{1}{3}$ . If the agent tries to move outside the grid or steps into a hole, the transition function renders it stationary (no change in state). In this case, the reward remains 0 unless the movement ends in the goal state, for which the agent receives a reward of 1. Paths blocked by holes or grid boundaries ensure that invalid transitions do not reward the agent and do not result in movement, maintaining the dynamics of the episode. The following sections detail the process behind optimizing a policy  $\pi$  for a given MDP and how this can be affected by a slippery or non-slippery Frozen Lake.

## II. Markov Decision Process Problem Definition

The Frozen Lake map used for this experiment is shown below. For this experiment, let us also state that `slippery = True` for completeness.

$$\begin{bmatrix} S & F & F & F \\ F & H & F & H \\ F & F & F & H \\ H & F & F & G \end{bmatrix}$$

### A. State Space

The state space  $S$  consists of all the positions in the Frozen Lake map grid. This is shown in row-wise format below:

$$S = \{0, 1, 2, \dots, 15\}$$

These states correspond to a grid index on the map, where each state (0-indexed) corresponds to a Start, Frozen, Hole, or Goal tile in a left-to-right, top-to-bottom order. Thus, the array ["SFFF", "FHFH", "FFFH", "HFFG"] maps exactly to the Frozen Lake configuration above.

### B. Action Space

The action space  $A$  is defined as:

$$A = \{0, 1, 2, 3\}$$

where 0 = Left, 1 = Down, 2 = Right, and 3 = Up.

### C. Transition Function of Probabilities

Let  $P(s' | s, a)$  denote the probability of transitioning from state  $s$  to state  $s'$  using action  $a$ . For terminal states (holes  $H$  or goal  $G$ ):

$$P(s' = s | s, a) = 1, \quad \forall a \in A$$

For non-terminal states ( $s \in \{0, 1, 2, 3, 4, 6, 8, 9, 10, 13, 14\}$ ):

$$\begin{cases} P(s' = \text{Intended Direction}) = \frac{1}{3} \\ P(s' = \text{Move Up}) = \frac{1}{3} \\ P(s' = \text{Move Down}) = \frac{1}{3} \end{cases}$$

#### D. Discount Factor

The MDP has a constant discount factor:

$$\gamma = 0.95$$

#### E. Reward Function

The reward function is:

$$R(s, a, s') = \begin{cases} 1 & \text{if } s' = 15 \text{ (goal)} \\ 0 & \text{otherwise} \end{cases}$$

### III. Methods

The main objective is to obtain the optimal value function  $V^*(\cdot)$  and the optimal policy  $\pi^*(\cdot)$  using three distinct model-free approaches: on-policy MC control, SARSA, and Q-learning. These model-free methods resolve the issues caused by continuous state and action spaces. Even when the number of states and actions is very large, the optimal value function and policy can be estimated through a data-driven approach. In the following discussions, each of the three model-free approaches is outlined.

#### A. Parameters

- $\epsilon$  (Epsilon-Greedy Parameter): Determines how often the agent explores by picking a random action instead of the current estimated best action.
- $\gamma$  (Discount Factor): Used to discount future expected returns (set to 0.95).
- `slippery`: If True, follow the defined MDP from Section II; otherwise, transitions are deterministic.

#### B. Pseudocode - First-Visit Monte Carlo Control: On-Policy

**On-policy first-visit MC control (for  $\epsilon$ -soft policies), estimates  $\pi \approx \pi_*$**

```

Algorithm parameter: small  $\epsilon > 0$ 
Initialize:
   $\pi \leftarrow$  an arbitrary  $\epsilon$ -soft policy
   $Q(s, a) \in \mathbb{R}$  (arbitrarily), for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ 
   $Returns(s, a) \leftarrow$  empty list, for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ 

Repeat forever (for each episode):
  Generate an episode following  $\pi$ :  $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$ 
   $G \leftarrow 0$ 
  Loop for each step of episode,  $t = T-1, T-2, \dots, 0$ :
     $G \leftarrow \gamma G + R_{t+1}$ 
    Unless the pair  $S_t, A_t$  appears in  $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$ :
      Append  $G$  to  $Returns(S_t, A_t)$ 
       $Q(S_t, A_t) \leftarrow \text{average}(Returns(S_t, A_t))$ 
       $A^* \leftarrow \text{argmax}_a Q(S_t, a)$  (with ties broken arbitrarily)
    For all  $a \in \mathcal{A}(S_t)$ :
       $\pi(a|S_t) \leftarrow \begin{cases} 1 - \epsilon + \epsilon/|\mathcal{A}(S_t)| & \text{if } a = A^* \\ \epsilon/|\mathcal{A}(S_t)| & \text{if } a \neq A^* \end{cases}$ 

```

**Fig. 1 Sutton and Barto 5.4: On-Policy First-Visit MC Control.**

### C. Pseudocode - SARSA: On-Policy, TD(0)

```
Sarsa (on-policy TD control) for estimating  $Q \approx q_*$ 

Algorithm parameters: step size  $\alpha \in (0, 1]$ , small  $\varepsilon > 0$ 
Initialize  $Q(s, a)$ , for all  $s \in \mathcal{S}^+$ ,  $a \in \mathcal{A}(s)$ , arbitrarily except that  $Q(\text{terminal}, \cdot) = 0$ 
Loop for each episode:
  Initialize  $S$ 
  Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\varepsilon$ -greedy)
  Loop for each step of episode:
    Take action  $A$ , observe  $R, S'$ 
    Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\varepsilon$ -greedy)
     $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)]$ 
     $S \leftarrow S'; A \leftarrow A'$ 
  until  $S$  is terminal
```

**Fig. 2 Sutton and Barto 6.4: On-Policy SARSA**

### D. Pseudocode - Q-Learning: Off-Policy, TD(0)

```
Q-learning (off-policy TD control) for estimating  $\pi \approx \pi_*$ 

Algorithm parameters: step size  $\alpha \in (0, 1]$ , small  $\varepsilon > 0$ 
Initialize  $Q(s, a)$ , for all  $s \in \mathcal{S}^+$ ,  $a \in \mathcal{A}(s)$ , arbitrarily except that  $Q(\text{terminal}, \cdot) = 0$ 
Loop for each episode:
  Initialize  $S$ 
  Loop for each step of episode:
    Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\varepsilon$ -greedy)
    Take action  $A$ , observe  $R, S'$ 
     $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$ 
     $S \leftarrow S'$ 
  until  $S$  is terminal
```

**Fig. 3 Sutton and Barto 6.5: Off-Policy Q-Learning**

### E. Code Implementation

The previously shown pseudocode algorithms were implemented using Python, the gymnasium library, the numpy library, and the matplotlib library for plotting. Additionally, the Frozen Lake environment (slippery or not) had to be unwrapped to access the transition and reward function. The policies that correspond to each agent are outputted as lists. The value functions are outputted as 4x16 matrices.

## IV. Results

### A. Mean Evaluation Return for Each Trained Agent

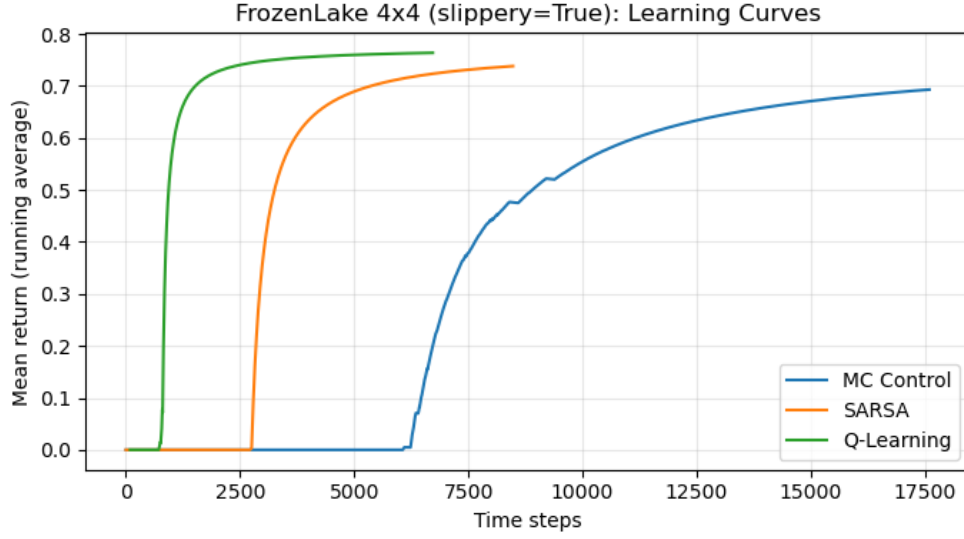


Fig. 4 Mean Evaluation Return - Slippy

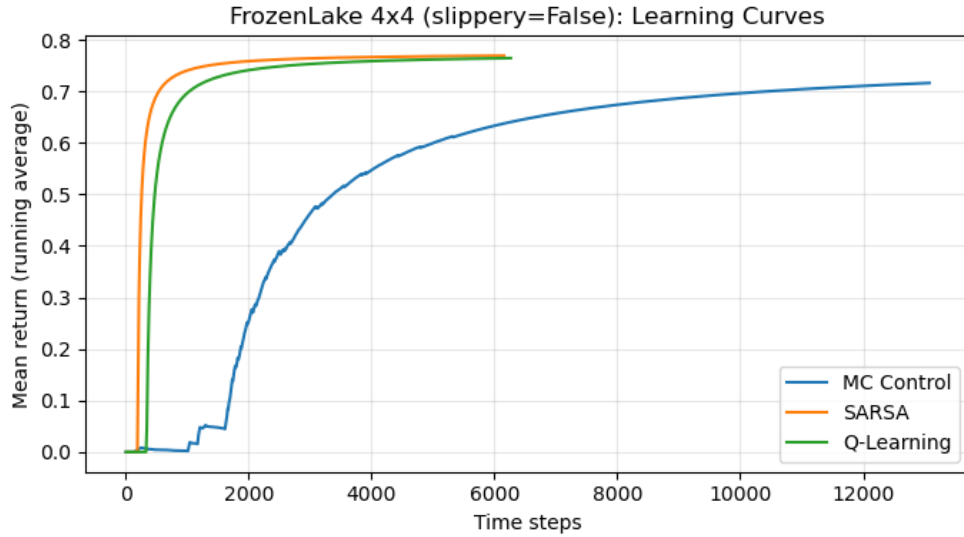


Fig. 5 Mean Evaluation Return - Not Slippy

In Figure 3, it can be observed that all three algorithms gradually improve their evaluation returns over time despite the inherent stochasticity of state transitions. Q-learning demonstrates faster learning and higher asymptotic performance than SARSA (orange) and Monte Carlo control (blue). This is because Q-learning is off-policy, and it can aggressively exploit optimal actions via the argmax function for the action-value function  $Q$  for its target. This leads to higher variance but faster convergence toward the optimal value. SARSA, being on-policy, updates using the next action actually taken under the  $\epsilon$ -greedy policy, which makes it more conservative and results in a smoother but slower learning. Monte Carlo control also exhibits the slower early learning due to its reliance on complete episodes before updates,

though it eventually approaches near optimal returns with many more time steps. Overall, Q-learning achieves the best balance between sample efficiency and asymptotic performance under the Frozen Lake’s slippery stochastic dynamics.

In Figure 4, all three algorithms exhibit faster and more stable learning compared to the slippery case. The deterministic transitions allow Q-learning and SARSA to converge rapidly to near-optimal performance, with Q-learning maintaining a slight edge in both speed and final return due to its efficient off-policy updates. Monte Carlo Control again demonstrates slower early learning, as it requires complete episodes to update its values, but eventually approaches optimality with sufficient time steps. Overall, the deterministic characteristics of the non-slippy environment enables all methods to learn more efficiently, with Q-learning achieving the best combination of sample efficiency and asymptotic performance, while SARSA closely follows.

## B. Policy and State Value Function Output for Each Trained Agent

Below are the policy and value function outputs for each trained agent.

```
SLIPPERY: True
MC Policy: [1 0 0 2 1 0 3 0 2 1 0 1 1 2 2 0]

Q Policy: [1 3 0 3 1 0 2 3 2 2 1 1 1 0 2 0]

SARSA Policy: [2 2 1 1 1 1 1 2 2 2 1 1 2 1 2 1]
```

Fig. 6 Policy Output for All Three Methods - Slippery

```
SLIPPERY: False
MC Policy: [1 0 0 3 1 0 3 1 2 1 0 2 0 2 2 2]

Q Policy: [1 0 3 3 1 0 2 1 2 2 1 2 0 2 2 1]

SARSA Policy: [1 1 0 0 1 0 1 2 2 2 1 1 0 3 2 0]
```

Fig. 7 Policy Output for All Three Methods - Not Slippery

```
MC Value Function:
[4.57541906e-01 7.73780937e-01 4.37868249e-01 2.39979576e-01]
[4.61274982e-01 9.58629311e-04 1.00680137e-04 9.38119569e-04]
[8.35199220e-02 7.87826855e-02 3.41246076e-02 8.0989136e-02]
[4.36311259e-01 2.86017604e-01 4.50964440e-01 3.04856388e-01]
[2.71448021e-01 8.14506250e-01 1.02028781e-03 2.66968132e-01]
[6.95414729e-01 5.28541848e-01 4.20386161e-01 4.96472024e-01]
[1.60262555e-01 1.04213812e-01 1.61019719e-01 1.73917513e-01]
[7.59026444e-01 2.80468162e-01 5.17256934e-01 3.25852765e-01]
[3.87753138e-01 5.60679834e-02 8.57375000e-01 3.15997384e-01]
[6.18830833e-01 9.02500000e-01 5.03618173e-01 5.58691303e-02]
[7.69452874e-01 4.13155420e-01 5.57297821e-01 4.72570114e-01]
[2.33318378e-01 6.61859245e-01 6.31662438e-01 6.09989098e-01]
[6.65324257e-01 9.81238796e-01 5.62939929e-01 6.37806262e-01]
[1.68449835e-01 0.30458504e-01 9.50000000e-01 7.57289446e-01]
[6.32467332e-01 7.42488427e-01 1.00000000e+00 5.34761883e-01]
[7.27604552e-01 4.89471811e-01 3.13655343e-01 5.08926837e-01]]

SARSA Value Function:
[[0.2574142 0.26141041 0.77378094 0.04752653]
[0.14117961 0.24556421 0.81450625 0.26298325]
[0.2178605 0.657375 0.15068051 0.27251922]
[0.5782667 0.70851457 0.6351293 0.23245277]
[0.26759763 0.28158201 0.26134288 0.26683724]
[0.48973187 0.71251138 0.04388428 0.43693061]
[0.27066963 0.9025 0.23023408 0.7997307]
[0.58872313 0.8018518 0.01604187 0.25539744]
[0.27115436 0.23324866 0.36055599 0.27406944]
[0.27740857 0.26711253 0.56453654 0.0738678 ]
[0.27769298 0.95 0.07115504 0.0312619]
[0.21503226 0.67203831 0.49536442 0.0308243]
[0.29633901 0.28852804 0.61999331 0.12023666]
[0.25541881 0.2792047 0.27740615 0.160812 ]
[0.27835768 0.27762442 1. 0.2803508 ]
[0.18132491 0.09614504 0.06673286 0.30153324]]

Q Value Function:
[[0.38053381 0.77378094 0.37598273 0.07846903]
[0.10254746 0.22784578 0.27313528 0.27536348]
[0.27407991 0.24661566 0.27271058 0.27201399]
[0.27228378 0.02449091 0.27024006 0.27393802]
[0.00405168 0.81450625 0.32030511 0.00192469]
[0.94657747 0.31157645 0.6827014 0.79710261]
[0.1106717 0.02507835 0.12029759 0.05716331]
[0.80806484 0.48832833 0.36518599 0.9931468 ]
[0.23275139 0.19320786 0.057375 0.22514943]
[0.16238334 0.18347483 0.9025 0.69349208]
[0.78916108 0.95 0.43367015 0.02584932]
[0.59031781 0.79342798 0.23071814 0.65191876]
[0.35672136 0.96425466 0.5584514 0.8509143 ]
[0.04580866 0.27155257 0.0554027 0.66831507]
[0.44150515 0.32253215 1. 0.64994906]
[0.9892063 0.67036846 0.1994473 0.84442758]]
```

Fig. 8 Value Function Output for All Three Methods - Slippery

```
MC Value Function:
[0.35084932 0.77378094 0.21029977 0.02608741]
[0.27298836 0.01186753 0.01342173 0.15073524]
[0.07749236 0.07304334 0.05439265 0.07141143]
[0.15668394 0.01991318 0.06475944 0.18262081]
[0.4082201 0.81450625 0.00029347 0.37322255]
[0.75676841 0.30417361 0.05943668 0.19485135]
[0.29408478 0.29717809 0.34548769 0.39483153]
[0.1386862 0.57656071 0.41103026 0.46812156]
[0.64353857 0.03381516 0.857375 0.31341357]
[0.63516849 0.9025 0.57794114 0.0597688 ]
[0.74125523 0.18814442 0.22247129 0.31910807]
[0.69855576 0.7134459 0.71040847 0.16472062]
[0.93999122 0.93250617 0.68798474 0.52755461]
[0.28431708 0.61219945 0.95 0.73822055]
[0.7551975 0.70129839 1. 0.64055498]
[0.18180995 0.59756375 0.71325769 0.05284581]]

SARSA Value Function:
[[0.17755383 0.77378094 0.08136185 0.08361533]
[0.39475362 0.08912682 0.56288499 0.52468364]
[0.8654551 0.1717157 0.22059699 0.0593376]
[0.54208007 0.12216638 0.42057575 0.09795267]
[0.18151564 0.81450625 0.1563857 0.1841025 ]
[0.97120751 0.64327216 0.87556331 0.67926308]
[0.1336855 0.16689605 0.12266908 0.1156938 ]
[0.28725985 0.85905602 0.09032207 0.47548779]
[0.185978 0.16238132 0.857375 0.18762621]
[0.15954135 0.12095294 0.9025 0.0217058 ]
[0.16576611 0.95 0.13276611 0.10542234]
[0.25712152 0.74702371 0.1825597 0.53680216]
[0.89504939 0.84580968 0.48866677 0.36396619]
[0.13782598 0.04063652 0.09669898 0.16864947]
[0.17837718 0.16880752 1. 0.16901072]
[0.75596583 0.4130171 0.04291126 0.21932081]]

Q Value Function:
[[2.094089717e-01 7.73780937e-01 1.88833226e-01 2.10277248e-01]
[2.58812632e-01 2.55542005e-01 5.02243229e-01 4.98107384e-01]
[3.35326146e-01 4.29986501e-01 8.7651710e-01 9.92461482e-01]
[5.25264977e-01 5.56406146e-01 1.45144444e-01 5.54117265e-01]
[1.09977816e-01 8.14506250e-01 9.06315506e-02 7.07362017e-02]
[2.88112709e-01 2.51823967e-01 3.75675196e-02 3.09860725e-02]
[0.99646102e-01 3.02008066e-01 9.30414622e-02 3.42135591e-02]
[1.12976627e-01 2.23566269e-01 1.56866706e-01 5.12263877e-01]
[6.56228441e-02 1.23298274e-01 8.57375000e-01 1.39732104e-01]
[6.85997165e-01 7.34640747e-01 9.02500000e-01 1.44155219e-01]
[3.46726839e-01 9.50000000e-01 5.84721354e-01 1.63177328e-01]
[2.41122522e-01 4.87028858e-01 9.25977328e-01 2.46354878e-02]
[7.69636990e-01 1.37958068e-01 7.02760034e-01 2.47237885e-01]
[2.77999608e-01 5.8024025e-01 7.1428215e-01 8.01108822e-04]
[5.68174444e-01 4.69276444e-01 1.00000000e+00 4.87286029e-01]
[8.79287761e-01 9.51338164e-01 1.05610232e-01 8.84823494e-02]]
```

Fig. 9 Value Function Output for All Three Methods - Not Slippery

## V. Conclusions

This report compared three RL methods: Monte Carlo, SARSA, and Q-learning algorithms in the 4×4 FrozenLake tabular environment under both stochastic (slippery) and deterministic (non-slippery) environments. Q-learning proved to be the most sample-efficient and consistently reached near-optimal performance fastest, especially when transitions were deterministic. SARSA showed smoother, more stable learning but converged more slowly—particularly in the presence of random transitions, which demonstrates its cautious, on-policy nature. Monte Carlo control improved steadily but was initially slower since it updates at the end of each episode, and with sufficient experience, it achieved returns close to the optimum. All algorithms performed better and learned faster in the deterministic (non-slippery) setting, which is expected. However, for small, discrete tasks like FrozenLake, Q-learning is best when prioritizing fast, optimal learning. Regarding the other two methods, SARSA is preferable when stability during training matters, and MC control is good for establishing a simple and reliable baseline given enough training data.

## References

- [1] Richard S. Sutton and Andrew G. Barto (2018). Reinforcement Learning: An Introduction, MIT Press