

# Reinforcement Learning with Frozen Lake

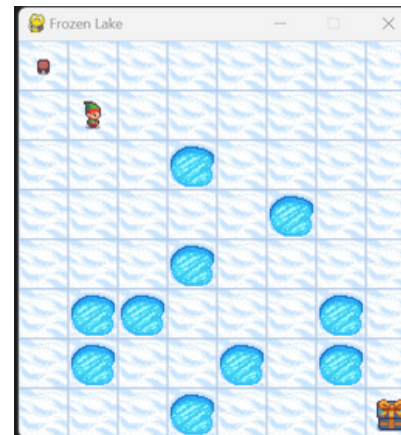
*José Antonio Sánchez González*

# Problem definition

In this Reinforcement Learning Framework, the algorithm involves crossing a frozen lake from Start(s) to Goal(G) without falling into any Holes(H) by walking over the Frozen(F) lake. The agent may not always move in the intended direction due to slippery nature of the frozen lake.

The agent takes a 1-element vector for actions. The observation is a value representing the agent's current position as  $\text{current\_row} * \text{nrows} + \text{current\_col}$  (it starts at 0). For example, the goal position in the 4x4 map can be calculated as follows:  $3 * 4 + 3 = 15$ .

The number of possible observations is dependent on the size of the map. For us is going to be an 8x8 map.



# Algorithms we use:

## Epsilon-Greedy Algorithm

eps=1

if random() < eps:

    select random action

else:

    select best action

#At the end of each episode:

eps = eps - decay rate

We are gonna work with **discrete data**, that is going to be stored in a file that is actually gonna be the q-table, to train the model.

## Q-Table

Q-Learning Formula

$$q[\text{state}, \text{action}] = q[\text{state}, \text{action}] + \text{learning\_rate} * (\text{reward} + \text{discount\_factor} * \max(q[\text{new\_state}, :]) - q[\text{state}, \text{action}])$$

# Frozen Lake

Frozen Lake is part of the Toy Text environments.

```
import gymnasium as gym  
gym.make("FrozenLake-v1")
```



Frozen lake involves crossing a frozen lake from Start(S) to Goal(G) without falling into any Holes(H) by walking over the Frozen(F) lake. The agent may not always move in the intended direction due to the slippery nature of the frozen lake.

## Action space:

- 0: LEFT
- 1: DOWN
- 2: RIGHT
- 3: UP

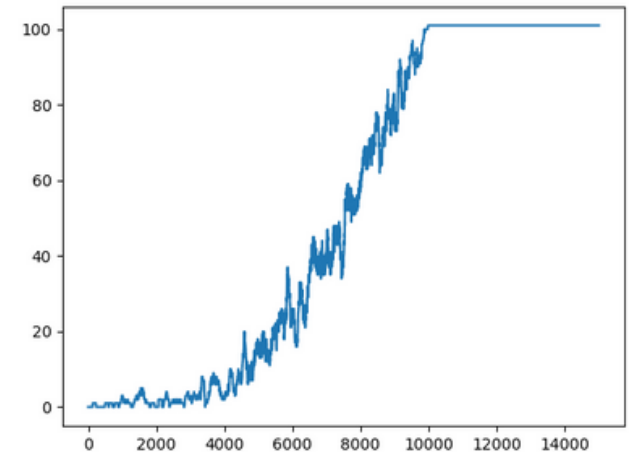
for more information read  
[gymnasium](#) documentation.

## FrozenLake: 15,000 episodes

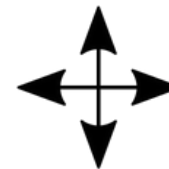
Y-axis for reward and X-axis for # of episodes.

we can see that at **100 episodes** is not nearly enough to learn. As the epsilon decreases, at **5000 episodes** it starts to use the Q-table to learn.

At **10,000 episodes**, epsilon is 0 this means it is only gonna follow the Q-table and it's not gonna learn anymore.

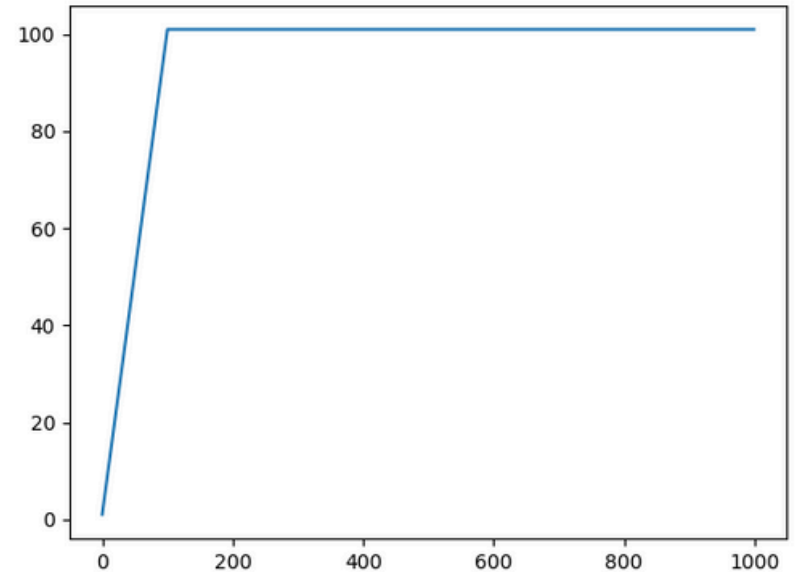


**after 10,000 episodes we get to the goal a 100% of the time!**



from 0 to 100 it's going to look like an increasing line.

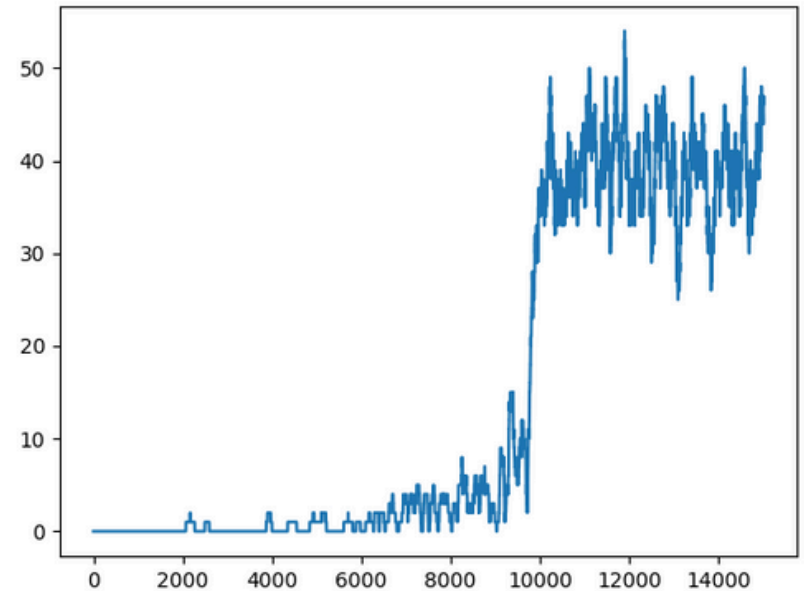
Then it gets 100% reward.



running for 1000 episodes

If we run 15,000 episodes and with slippery factor we get this result.

barely after 8,000 episodes it starts getting rewards.



**the slippery factor makes training a lot more difficult!**

Here is the full code if you want to peek at it.

```
GymRL.py X frozen_lake8x8.png
GymRL.py > ...
1 import gymnasium as gym
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import pickle
5
6 def run(epochs, is_training=True, render=False):
7
8     env = gym.make("FrozenLake-v1", map_name="8x8", is_slippery=False, render_mode="human" if render else None)
9
10    if(is_training):
11        q = np.zeros((env.observation_space.n, env.action_space.n)) #init a 64x4 array
12    else:
13        f = open('frozen_lake8x8.pkl', 'rb')
14        q = pickle.load(f)
15        f.close()
16
17    learning_rate_a = 0.9 #alpha or learnig rate
18    discount_factor_g = 0.9 #gamma or discount factor
19
20    epsilon = 1 # 1-100% random actions
21    epsilon_decay_rate = 0.0001 # epsilon decay rate
22    rng = np.random.default_rng() # random number generator
23
24    rewards_per_episode = np.zeros(epochs)
25
26    for i in range(epochs):
27        state = env.reset()[0] #states: 0 to 63, 0=top left corner, 63=bottom right corner
28        terminated = False #True hwnw fall in hole or reached goal
29        truncated = False #True when actions > 200
30
31        while(not terminated and not truncated):
32            if is_training and rng.random() < epsilon:
33                action = env.action_space.sample() # actions: 0=left, 1=down, 2=right, 3=up
34            else:
35                action = np.argmax(q[state,:])
36
37            new_state, reward, terminated, truncated, _ = env.step(action)
38
39            if is_training:
40                q[state,action] = q[state,action] + learning_rate_a * (
41                    reward + discount_factor_g * np.max(q[new_state,:]) - q[state,action]
42                )
43
44            state = new_state
45
46        epsilon = max(epsilon - epsilon_decay_rate, 0)
47
48        if(epsilon==0):
49            learning_rate_a = 0.0001
50
51        if reward == 1:
52            rewards_per_episode[i] = 1
53
54    env.close()
55
56    sum_rewards = np.zeros(epochs)
57    for t in range(epochs):
58        sum_rewards[t] = np.sum(rewards_per_episode[max(0, t-100):(t+1)])
59    plt.plot(sum_rewards)
60    plt.savefig('frozen_lake8x8.png')
61
62    if is_training:
63        f = open('frozen_lake8x8.pkl', "wb")
64        pickle.dump(q, f)
65        f.close()
66
67
68 if __name__ == '__main__':
69     arun(15000)
70
71     run(1000, is_training=True, render=False)
```

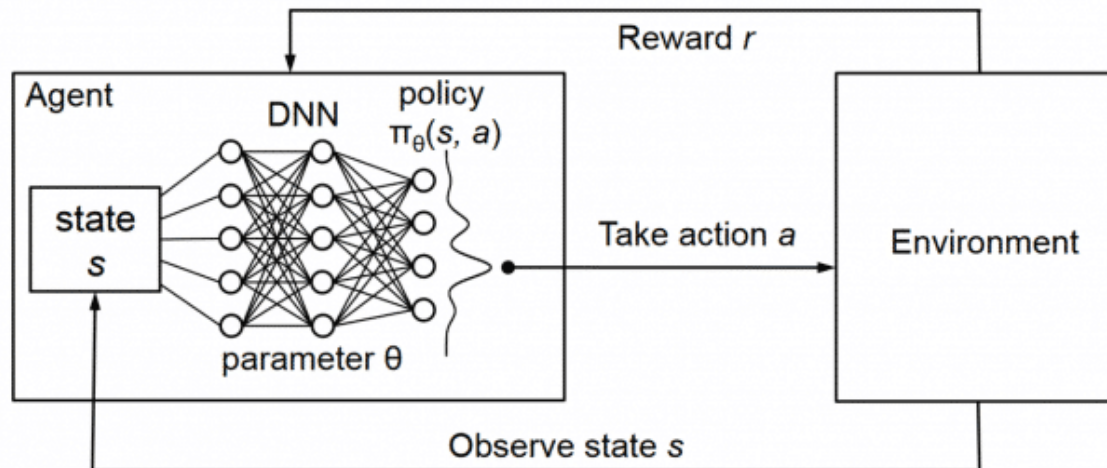


# Conclusions

- After we trained and tested the data there are instances where the goal is not met, it is highly dependant of the # of episodes and if you activate the slippery factor.
- The model can improve by adjusting the # of episodes and the epsilon factor
- There is an option to create a custom environment, therefore creating a custom environment might be a better option
- we can also use a **Deep Q-Learning** Formula which is the most promising as a custom environment.

## DQL Formula

$q[\text{state}, \text{action}] = \text{reward}$  if new\_state is terminal else  
 $\text{reward} + \text{discount\_factor} * \max(q[\text{new\_state}, :])$



**A DQL model would help us fine-tune the model and train it on different map sizes and worry less about the slippery factor, so it is truly advised to work on this further.**