

# Frozen Lake Assignment

Setup 1:

**FrequentRewards = False**

**T = 200**

- The Q table? why is it full of zeros?

In such setup we define lot of steps in episode and original Frozen Lake Reward. This lead us to the problem that agent almost never gets to the finish point to get Reward (100000 episodes are not enough to reach finish at least once) so never update Q-table (update with zeros)

- Is your agent learning anything?

Agent isn't learning because we are not passing any reward to the start of the chain, updating Q function  $Q(s, a)$  for each possible state and action based on formula

$$Q_{s,a} = \alpha * (R_{t+1} + \gamma \hat{Q}(S_2, \operatorname{argmax} Q(S_2, a))) - Q_{s,a} \quad (1)$$

Since Reward is always 0 (except last point which is non explored in most cases) and all Q-table initialized with zeros.

- What What happens when is\_slippery = False, versus when is\_slippery = True

In this setup **slippery** flag doesn't effect on result because even without it agent couldn't find path to finish. Obviously the constraint that there exist uncertainty when performing action should make it harder for agent to get optimal results.

Setup 2:

**FrequentRewards = True**

**Exploration\_rate = 0.01**

**Is\_slippery = True**

- How has your Q-table changed and why?

Now Q-table actually update because with some values. Agent receives reward all the time. So then we have non-zero values for function (1) to work with.

- Is your agent learning anything?

Agent gets penalty when falls into water, so with this constraint, to maximize reward gain, agent learn not to go there. Also other constraint prevents agent from going around by decreasing reward by a tiny rate each step. Finally agent still get positive reward when gets to finish. This happens more frequent because now agent avoid bumping into water each episode, this increase explored area, and probability of getting to finish is enough to get meaningful probabilities in Q-table.

This setup ensures that agent will learn to go to the finish, minimize number of steps and avoid falling into water.

Now pick a large value for exploration rate 0.1 or 0.5, how does this impact the total reward curve?

Large exploration rate decrease a success rate of agent significantly. This is because we decide to make random decision in 50% steps, so we rather explore all possible options, instead of focusing and optimizing actual best path.

It is possible to solve this by using some epsilon param which dynamically represent exploration\_rate. It is initialized as 1 and slowly decreased to 0 during learning. This ensure that firstly agent explore and for Q-table, and gradually going to master optimal path by following Q-table up to 100% of steps.

Now set a very short episode length say  $T = 20$

Such hyperparameter is too small for the task. It is almost impossible to find the finish by random, so agent never gets to the finish, without understanding goal(positive reward) agent will never form the path to optimize, even with such path exist.

In our task 20 steps is too small ammount, taking into account that minimal path is 16, but slippery makes probability extremely small  $(1 - P(\text{slip}))^{16}$

Reflect on why having frequent rewards has such an impact on learning

The update of the Q-table is more dynamic and consiquent. As the agent interacts with the environment, it collects rewards at many steps, not just at the goal. Each action taken by the agent leads to some sort of immediate feedbcack, which helps in faster learning and adjustment of strategies to maximize rewards. This frequent feedback leads to numerous updates in the Q-table, reflecting learning from various states and actions. Each entry in the Q-table represents the expected future rewards for a given state-action pair, and with more frequent rewards, these values become more refined all the time.

Note that here each experience is used only once to update the Q-table, what is the disadvantage of this for rare events?

Each experience tuple influences the Q-table only once, at the time it occurs in episode. So for rare events, this can mean too few updates to capture this goal or penalty by agent, potentially leading to underestimation of important but rare strategies.

Do you have a better solution?

I researched the solution which is commonly used with such problem

**Experience Replay** is a replay memory technique used in reinforcement learning where we store the agent's experiences at each time-step, in a data-set , pooled over many episodes into a replay memory. We then usually sample the memory randomly for a minibatch of experience, and use this to learn off-policy [\[link\]](#)

At any point in learning, the agent remembered the best episode up to that point. After every four episodes, it replayed this remembered episode, learning from it as if it were a new episode. [Reinforcement Learning: An Introduction, Richard S. Sutton and Andrew G. Barto]

Such technique capture the best result and repeat them, so good episodes have more effect on Q-table by repeating them. This solve the issue that rare but significant example have small impact on learning process.

What are other methods to explore other than epsilon-greedy, do you have any proposals?

## Classic Exploration Strategies

- **Epsilon-greedy:** The agent does random exploration occasionally with probability  $\epsilon$  and takes the optimal action with probability  $1 - \epsilon$ .
- **Upper confidence bounds:** The agent selects the greediest action to maximize the upper confidence bound  $\hat{Q}_t(a) + \hat{U}_t(a)$ , where  $\hat{Q}_t(a)$  is the average rewards associated with action  $a$  up to time  $t$  and  $\hat{U}_t(a)$  is a function reversely proportional to how many times action  $a$  has been taken.
- **Boltzmann exploration:** The agent draws actions from a boltzmann distribution (softmax) over the learned Q values, regulated by a temperature parameter  $\tau$ .
- **Thompson sampling:** The agent keeps track of a belief over the probability of optimal actions and samples from this distribution.

The following strategies could be used for better exploration in deep RL training when neural networks are used for function approximation:

- **Entropy loss term:** Add an entropy term  $H(\pi(a | s))$  into the loss function, encouraging the policy to take diverse actions.
- **Noise-based Exploration:** Add noise into the observation, action or even parameter space [\[link\]](#)

In this document I wont provide detailed explanation of this strategies