

ECE 1508: Reinforcement Learning

Chapter 1: Introduction

Ali Bereyhi

`ali.bereyhi@utoronto.ca`

Department of Electrical and Computer Engineering
University of Toronto

Fall 2025

Frozen Lake Game



The guy needs to get to the treasure

- It could walk over frozen cells
- If it runs to a shallow cell it falls into water

If he decides for a direction, it could

- move to direction with probability $1 - p$
- slip with probability p to **another** direction

*It wants to learn the **best way**, i.e.,*

with minimal chance of failing

Playing in RL Framework

For this problem, we should first define *the main three components*

- *State*
- *Action*
- *Reward*

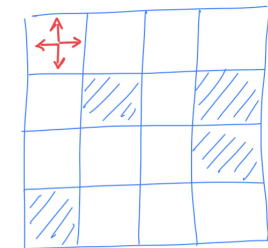
State \equiv Location of the guy on the game board

0	1	2	3
4	5	6	7
8	9	10	11
12	13	14	15

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

Playing in RL Framework

Action \equiv Possible moves on the board

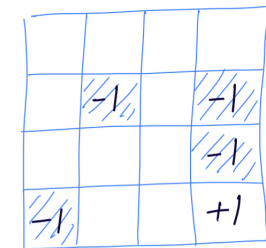


We can represent them with some integer

$$\mathbb{A} = \{0 \equiv \text{left}, 1 \equiv \text{down}, 2 \equiv \text{right}, 3 \equiv \text{up}\}$$

Playing in RL Framework

Reward \equiv Possible outcomes of the game



	-1		-1
			-1
-1			+1

We have three possible cases

$$R \in \{0 \equiv \text{ongoing}, 1 \equiv \text{winning}, -1 \equiv \text{losing}\}$$

Playing in RL Framework

- + Can we describe the *environment* in this problem?
- Sure!

Rewarding Function

Reward of each *state-action* pair is a random variable

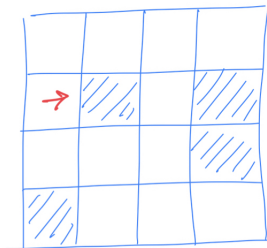
$$\mathcal{R}(s, a) = \begin{cases} \text{intended reward} & 1 - p \\ \text{reward of slipping direction} & p \end{cases}$$

Transition Function

Transition of each *state-action* pair is also a random variable

$$\mathcal{P}(s, a) = \begin{cases} \text{intended state} & 1 - p \\ \text{state of slipping direction} & p \end{cases}$$

Playing in RL Framework

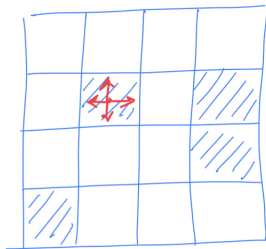


For example, in the above image we could say

$$\mathcal{P}(\textcolor{green}{4}, \textcolor{red}{2}) = \begin{cases} 5 & 1-p \\ 0 \text{ or } 8 & p \end{cases} \quad \text{and} \quad \mathcal{R}(\textcolor{green}{4}, \textcolor{red}{2}) = \begin{cases} -1 & 1-p \\ 0 & p \end{cases}$$

Terminal State

- + But, how can we specify the end of the game?
- We need to define a **terminal state**



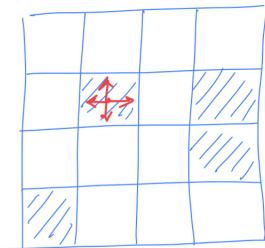
We can think of the end as a state which does not change anymore

$$\mathcal{P}(\mathbf{5}, a) = \mathbf{5}$$

for all **actions** $a \in \mathbb{A}$

Terminal State

- + What will be the **reward** then?
- We keep getting **zero reward**



We can think of the end as a state which does not change anymore

$$\mathcal{R}(\textcolor{green}{5}, \textcolor{red}{a}) = 0$$

for all **actions** $a \in \mathbb{A}$

Terminal State

Terminal State

State s is called terminal if for any action $a \in \mathbb{A}$, we have

$$\mathcal{P}(s, a) = s \quad \text{and} \quad \mathcal{R}(s, a) = 0$$

A basic property of **terminal state** is that it has **zero value**: say s is terminal

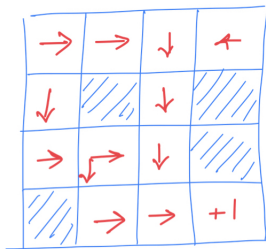
$$\begin{aligned} v_{\pi}(s) &= \mathbb{E}_{\pi} \{G_t | S_t = s\} = \mathbb{E}_{\pi} \{R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | S_t = s\} \\ &= \mathbb{E}_{\pi} \{0 + \gamma 0 + \gamma^2 0 + \dots | S_t = s\} = 0 \end{aligned}$$

Episode

*The trajectory from starting state to a terminal state is often called an **episode***

Playing in RL Framework: *Policy*

Policy in this problem describes the *planned path towards the goal*

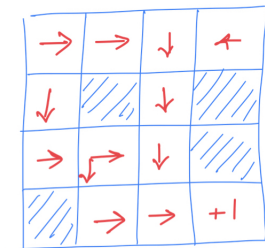


For instance in the above policy we have

$$\pi(a|1) = \begin{cases} 1 & a = 2 \\ 0 & a \neq 2 \end{cases} \quad \text{and} \quad \pi(a|9) = \begin{cases} 0.5 & a = 2 \\ 0.5 & a = 1 \\ 0 & a = 0, 3 \end{cases}$$

Playing in RL Framework: *Value*

Value describes the *chance of getting to the goal if we play the given policy*

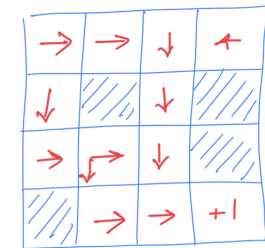


For instance, if slipping probability is $p = 0$ the above policy leads to

$$v_{\pi}(s) = \begin{cases} 1 & s \text{ is not terminal} \\ 0 & s \text{ is terminal} \end{cases}$$

Playing in RL Framework: *Optimal Policy*

Policy that **maximizes** the chance of getting to the goal



For instance, if slipping probability is $p = 0$ the above policy is **optimal**

Gymnasium



*Gymnasium is an API with some pre-implemented **environments***

- *We can call environments and access the **state** and **reward***
- *We can interact through **actions***
- *This can be a powerful toolkit to test our RL algorithms*
 - ↳ *At some point, we may prefer implementing our own environment*

Gymnasium: Trying Frozen Lake Game

We can call the frozen lake environment easily in Gymnasium

```
1  import gymnasium as gym
2
3  def test():
4      # Create FrozenLake instance
5      env = gym.make('FrozenLake-v1', map_name="4x4", is_slippery=False, render_mode='human')
6
7      env.reset() # Initialize to state 0
8      terminated = False # True when agent falls in hole or reached goal
9      truncated = False # True when agent takes more than 200 actions
10
11     while(not terminated and not truncated):
12
13         action = env.action_space.sample()
14
15         # Execute action
16         state, reward, terminated, truncated, _ = env.step(action)
17         print(reward)
18         print(state)
19     env.close()
```

We will play around with it in Assignment 1