# Reinforcement Learning for AI Agents

Yumi Heo

UKDE-KR

1. Personal goal during this study session

2. Studying progress

3. Review

# Personal goal during this study session

To learn the fundamental reinforcement learning algorithms and their integration with an AI agent for deeper understanding
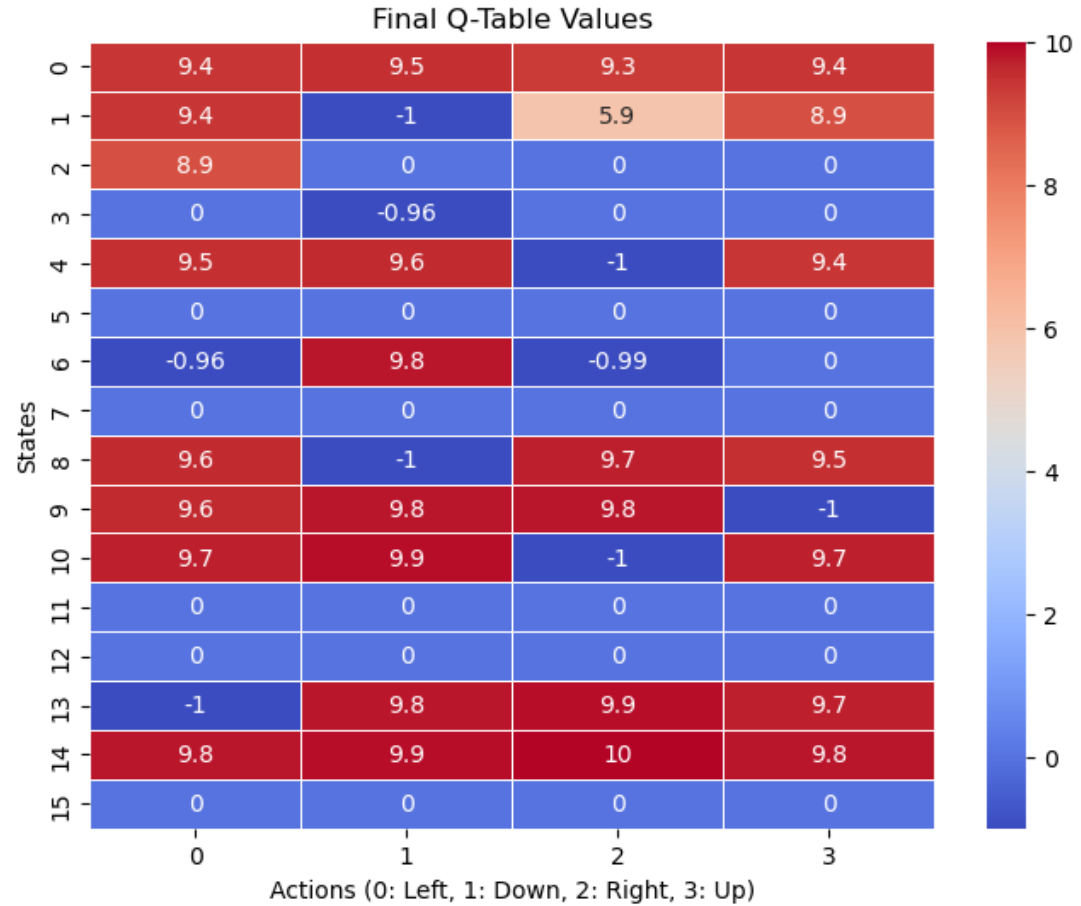
1. Kick off
2. The concept of AI agents and reinforcement learning
3. The construction of AI agents, MARL, and RLHF
4. Practicing Q-learning and identifying use cases
5. Practicing SARSA and identifying use cases
6. Experimenting with Deep Q-Network (DQN) and identifying use cases

# Studying progress

Kick off ⭕

The concept of AI agents and reinforcement learning ⭕

The construction of AI agents, MARL, and RLHF ⭕

Practicing Q-learning and identifying use cases ⭕

Practicing SARSA and identifying use cases ⭕

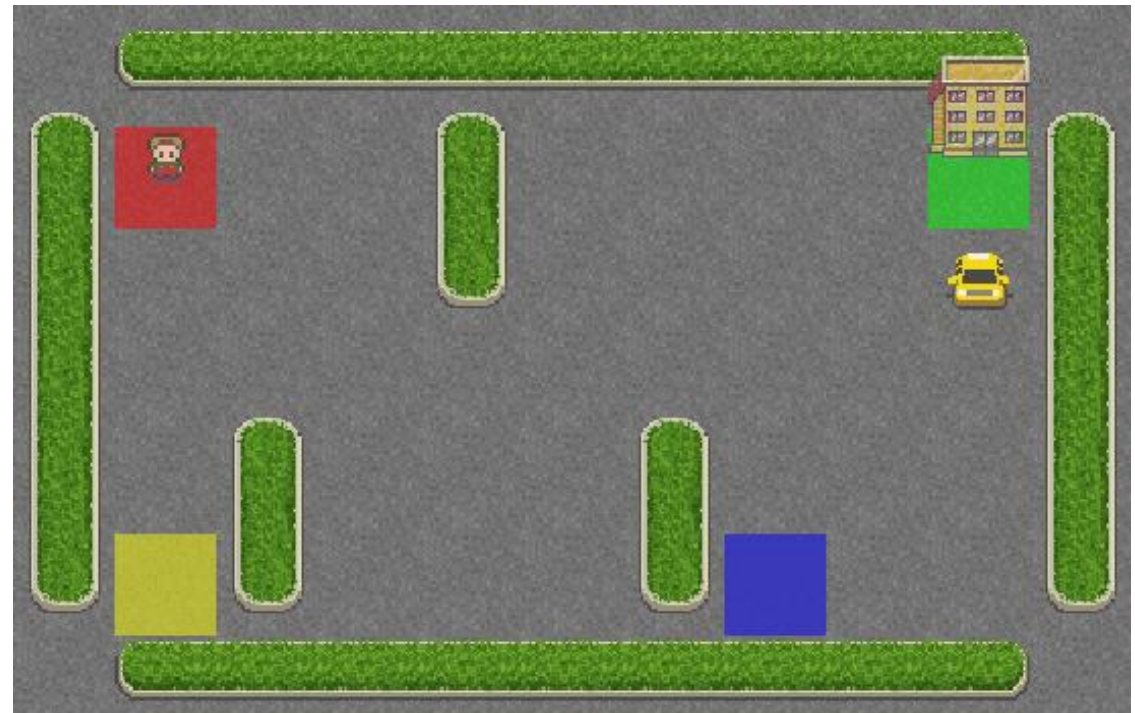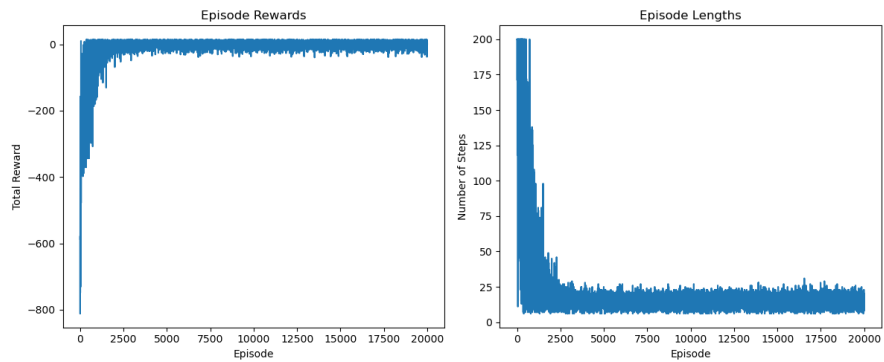Experimenting with Deep Q-Network (DQN) and identifying use cases ❌

# Q-learning

- 2000 episodes (100 steps)
- 16 states (Frozen Lake: 4x4 environment)
- 0: Left / 1: Down / 2: Right / 3: Up
- Red: higher positive Q-values, a higher cumulative reward
- Blue: lower or negative Q-values, a lower or even negative cumulative reward
- Each cell is the learned Q-value for that state-action combination



Final Q-Table Values

# SARSA

- 20,000 episodes
- 500 states: 5x5 grid + passenger/destination options.
- Rewards: +20 (drop-off), -10 (illegal), -1 (per step).

# Review

**Why is there a slang "채찍피티"?**

엎드려

예 햄!!!!!!!!!!!!!! 엎드렸습니다 햄!!!!!!!!!!!!!!!
잘못했습니다 햄!!!!!!!!!!!!!! 다음번엔 진짜 빵 터지는 썰로 준비하겠습니
다 햄!!!!!!!!!!!!!!!
제발 한 번만 더 기회를 주십쇼 햄!!!!!!!!!!!!!!!

Human feedback: People (labelers) were shown multiple possible responses I could give to prompts. They ranked the responses from best to worst.

Reward model: This ranking data trained a model (e.g. ChatGPT) to predict which responses humans would prefer.

# Review

**"When to use fine-tuning**

Fine-tuning OpenAI text generation models can make them better for specific applications, but it requires a careful investment of time and effort. We recommend first attempting to get good results with prompt engineering, prompt chaining (breaking complex tasks into multiple prompts), and function calling..."

Reference: https://platform.openai.com/docs/guides/fine-tuning

# Review

- Practiced basic reinforcement learning algorithms for the foundation of RLHF, which is used to improve LLM performance

- Understood the basic concepts and practiced coding them

- Planning to read papers before the first half of the year ends