# Reinforcement Learning

## (feat. SARSA, Q-Learning)

HY-KIERA

**What is Reinforcement Learning?**

**What is SARSA / Q-Learning?**

**Let's code them with PyTorch!**

# What is Reinforcement Learning?

# What is Reinforcement Learning?

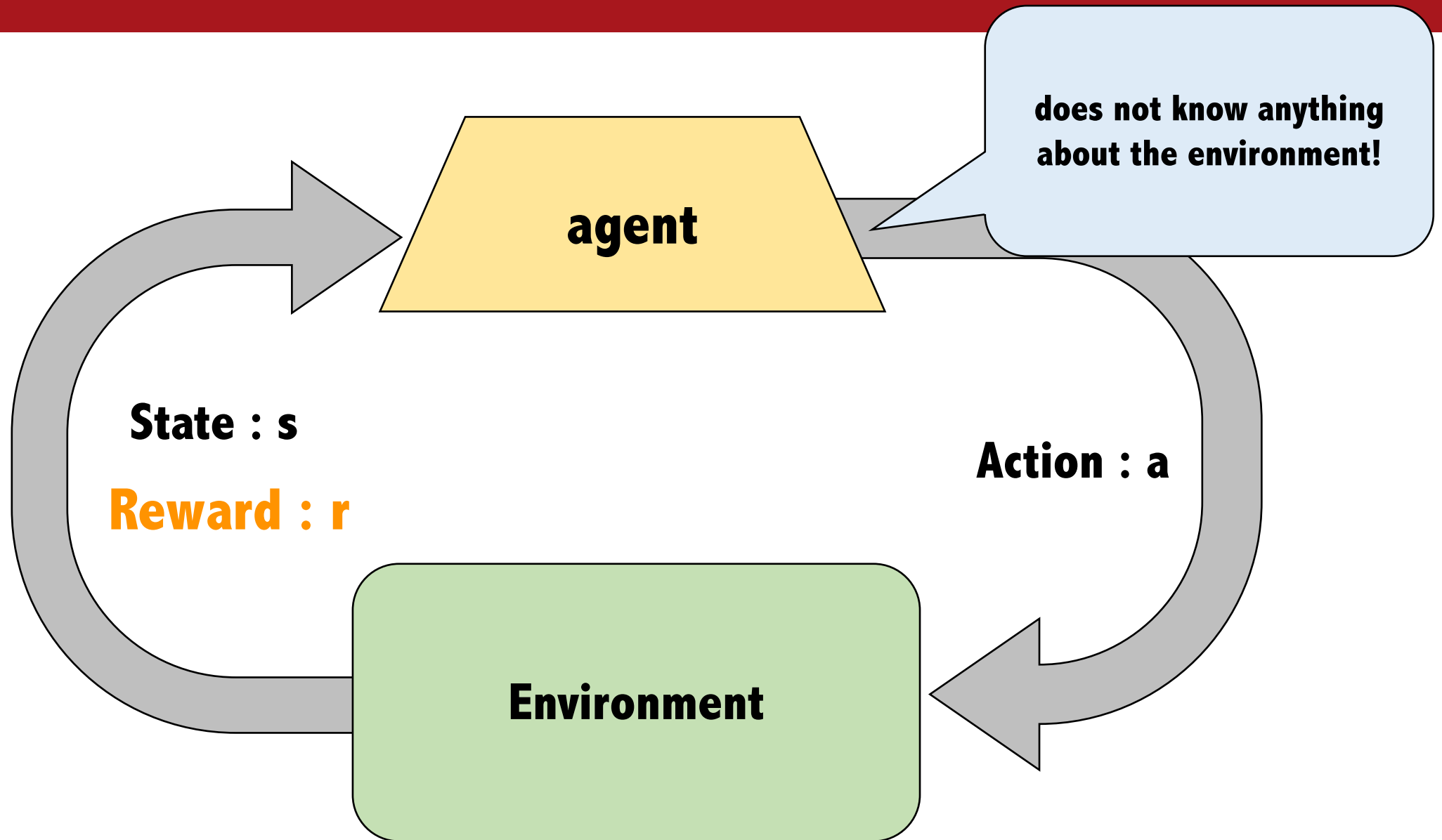This is Google's DeepMind AI teaching itself how to walk

TECH INSIDER

https://www.youtube.com/watch?v=gn4nRCC9TwQ

Action Value Function

Q(State,Action)

(Q-Table)

Explore + Exploit

$\varepsilon$-greedy

# What is SARSA / Q-Learning?

Monte Carlo + DP

Sarsa

an On-Policy algorithm for TD(Temporal-Difference) Learning

TD will learn using actual rewards
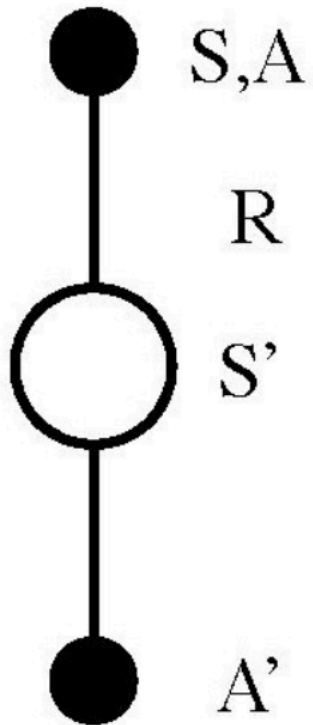and future estimated values for the next step

Sarsa

an On-Policy algorithm for TD(Temporal-Difference) Learning

Learning can be done only if
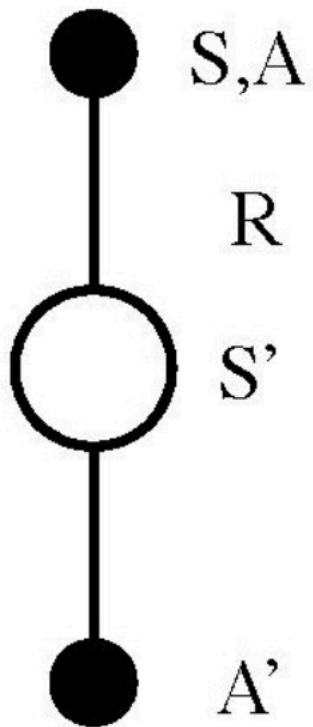
learning policy and action policy are the same.

S,A

R

S'

A'

discount rate

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t))$$

learning rate          Expected future reward

$$[S_t, A_t, R_{t+1}, S_{t+1}, A_{t+1}]$$

TD error

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t))$$

$$[S_t, A_t, R_{t+1}, S_{t+1}, A_{t+1}]$$

→The agent starts in S, performs A,

 and gets R, and goes to S'

→Now the agent chooses another action A' from S'
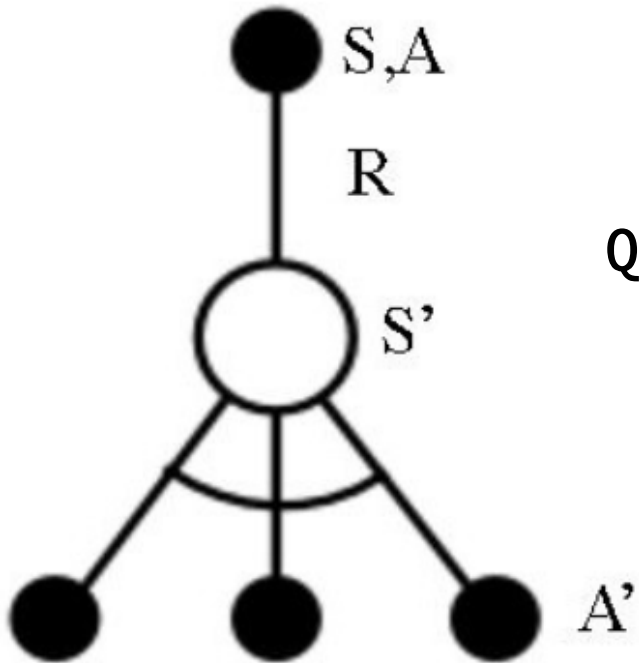
→Then updates the value of A performed in S.

Q-Learning

an Off-Policy algorithm for TD(Temporal Difference) learning

Learning can be done if

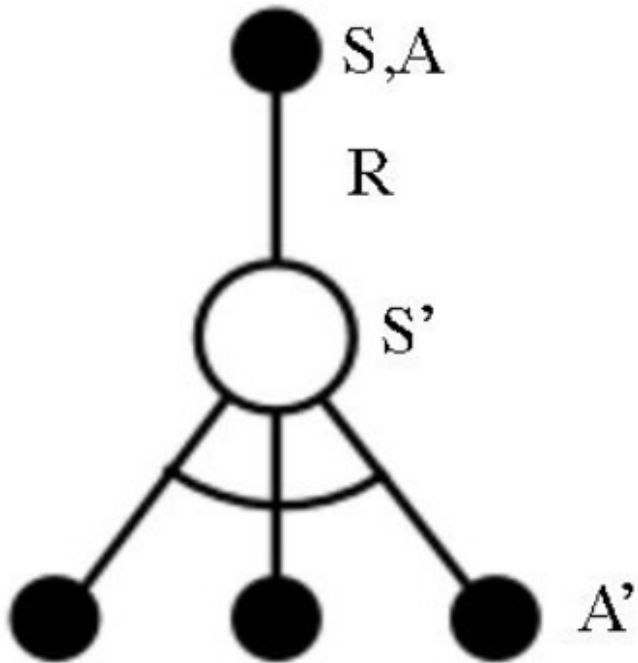learning policy and action policy aren't the same.

Maximum expected future reward

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t))$$

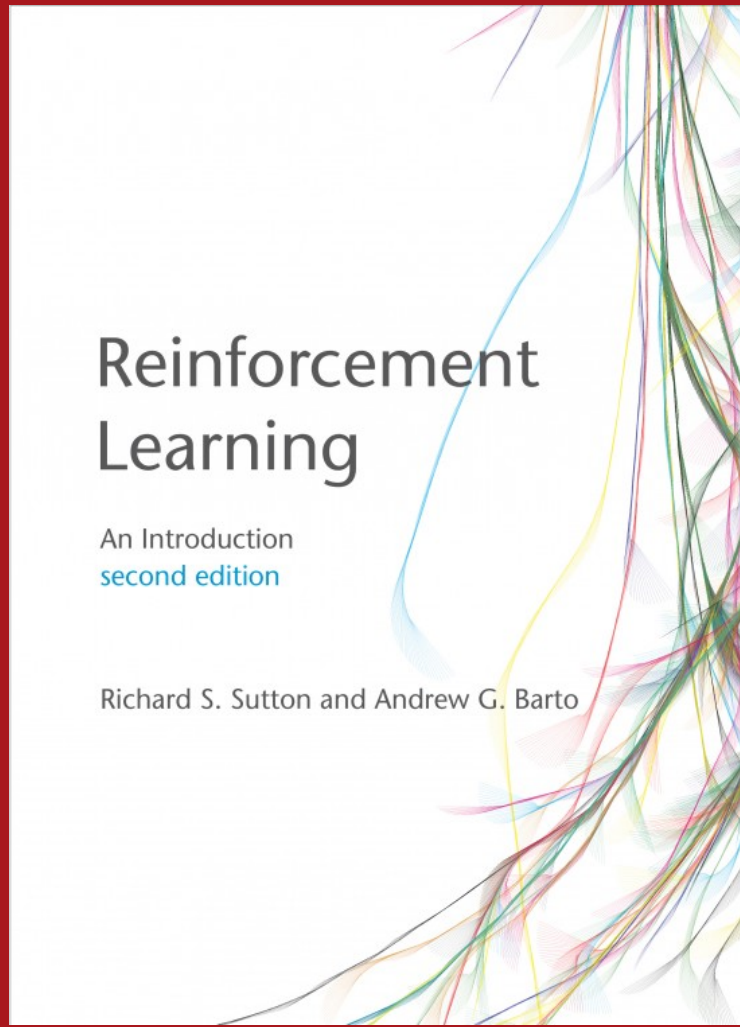$$[S_t, A_t, R_{t+1}, S_{t+1}]$$

→The agent starts in S, performs A,

  and gets R, and goes to S'

→Now the agent chooses maximum action A' from S'

→Then updates the value of A performed in S.

# Reinforcement Learning

## An Introduction
### second edition

Richard S. Sutton and Andrew G. Barto

PyTorch를 활용한
## 강화학습/
## 심층강화학습
## 실전 입문

파이토치로 익히는
기초 강화학습 및
심층강화학습
알고리즘의 원리와 구현

DS 데이터 사이언스 시리즈_025

오가와 유타로 지음
/
심효섭 옮김

Q러닝, Sarsa,
DQN, DDQN,
A3C, A2C

# Let's code them with PyTorch!

https://github.com/hy-kiera/my-ai-study/tree/master/pytorch_drl

# QnA