

# 강화학습과 최적화

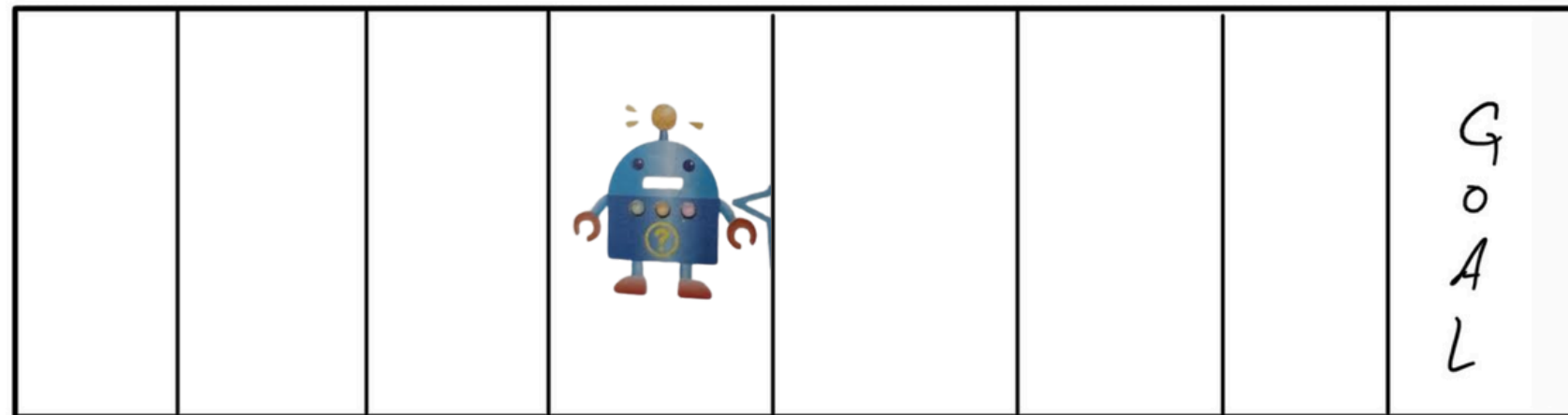
# 목차

---

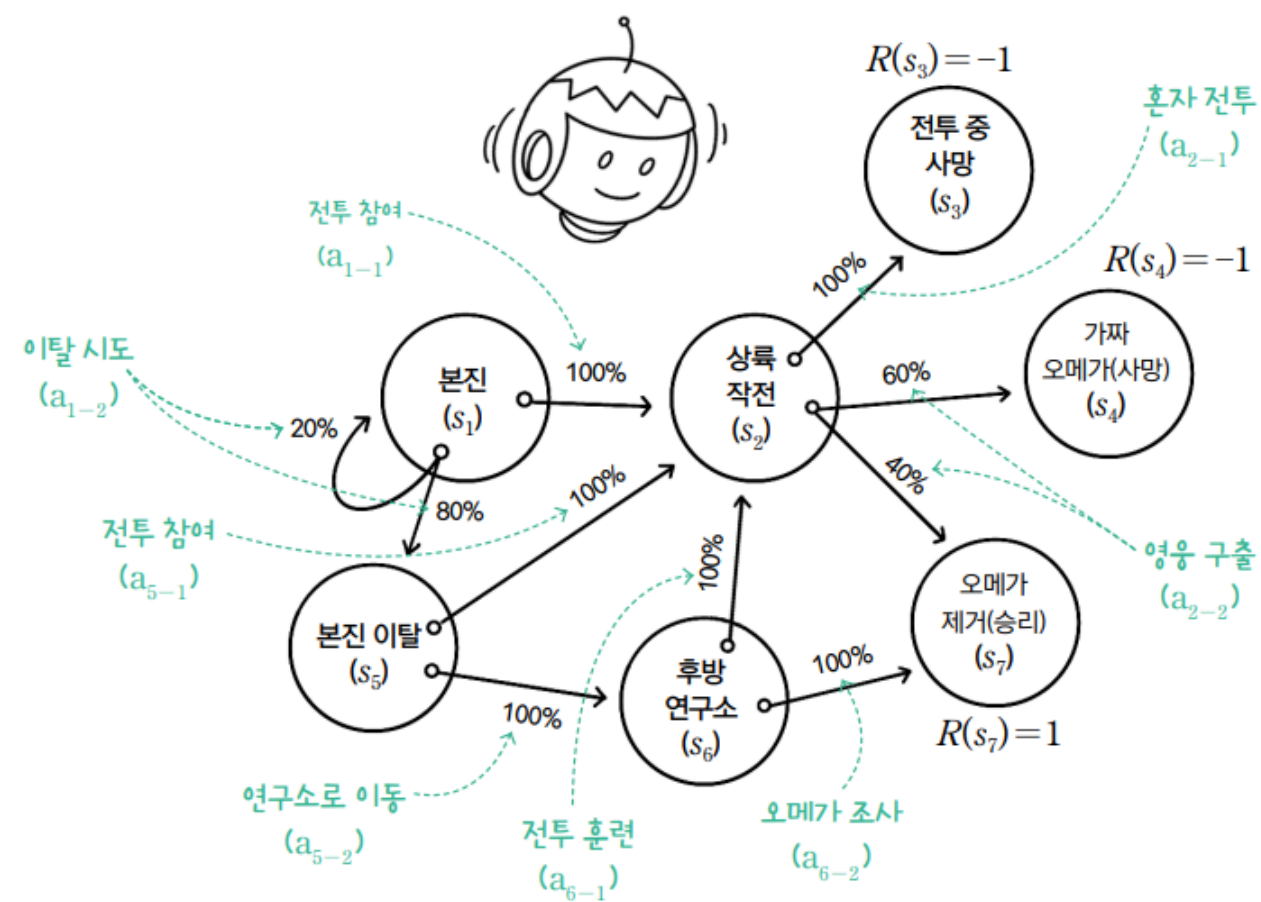
1. 강화학습이란?
2. 마르코프 결정 과정(MDP)
3. 가치 기반 강화학습
4. 정책 기반 강화학습
5. 정책 기반 강화학습을 적용한 하이퍼파라미터 최적화

# 강화학습이란?

머신러닝의 한 분야로, 의사 결정에 주로 사용되는 방법



# Markov Decision Process



$$p(a_1 | s_0 a_1 s_1)$$

$$p(a_1 | s_1)$$

	<h1>강화학습의 종류</h1>	
--	-------------------	--

Value-based

가치 기반

Policy-based

정책 기반

	<h1>Value-based</h1>	
--	----------------------	--

목적 함수로 상태 가치 함수, 행동 가치 함수를 설정. 보통 행동 가치 함수를 최대화하는 정책이 최적 정책이다.

상태 가치 함수

$$V(s) = E\left[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t) | s_0 = s\right]$$

	<h1>Value-based</h1>	
--	----------------------	--

목적 함수로 상태 가치 함수, 행동 가치 함수를 설정. 보통 행동 가치 함수를 최대화하는 정책이 최적 정책이다.

행동 가치 함수

$$Q(s, a) = E\left[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t) | s_0 = s, a_0 = a\right]$$

# Bellman-equation

상태 가치 함수, 행동 가치 함수를 재귀적으로 정의

상태 가치 함수  $V(s) = E[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t) | s_0 = s]$



# Q-learning

대표적인 가치 기반 강화학습으로, Q(행동 가치 함수)를 최대화하는 방향으로 학습한다.


# Policy-based

가치 기반의 문제

- 1) 연속적인 문제에 적용하기 힘들
- 2) 상황에 대한 행동이 한 개로 결정

Die		Goal		Die

	<h1>Policy-based</h1>	
--	-----------------------	--

정책 기반의 특징

- 1) 연속적인 문제에 적용 가능
- 2) 상황에 대한 행동이 평균과 분산으로 결정

	<h1>Policy-based</h1>	
--	-----------------------	--

보상 자체를 목적함수로 두어 사용

---

$$J \equiv E[G_0] = \int_{s_0, a_0, \dots} G_0 P(s_0, a_0, \dots) ds_0, a_0, \dots$$

---

	<h1>Optimization</h1>	
<p>Optuna</p> <p>C : 4.77</p> <p><math>\gamma</math> : 0.21</p> <p>소요 시간 : 3초</p> <p>스코어 : 0.9867</p>	<p>정책 기반</p> <p>C : 20.01</p> <p><math>\gamma</math> : 0.01</p> <p>소요 시간 : 18분 21초</p> <p>스코어 : 0.9733</p>	

---

# 느낀 점

강화학습의 개념과 종류에 대해 알게 되었고, 단순한 최적화 문제에서는 강화학습보다 optuna와 같은 비교적 간단한 방법이 더욱 효율적일 수 있다는 사실을 알게 되었다.

---

---

**Thank you.**

---