# TUTORIAL #11

Session 11
**Reinforcement Learning**

# LECTURE OVERVIEW

**01** | **Intro to RL**
What is it?
Why it is so important?

**02** | **RL Math**
Understanding of the
fundamental  math
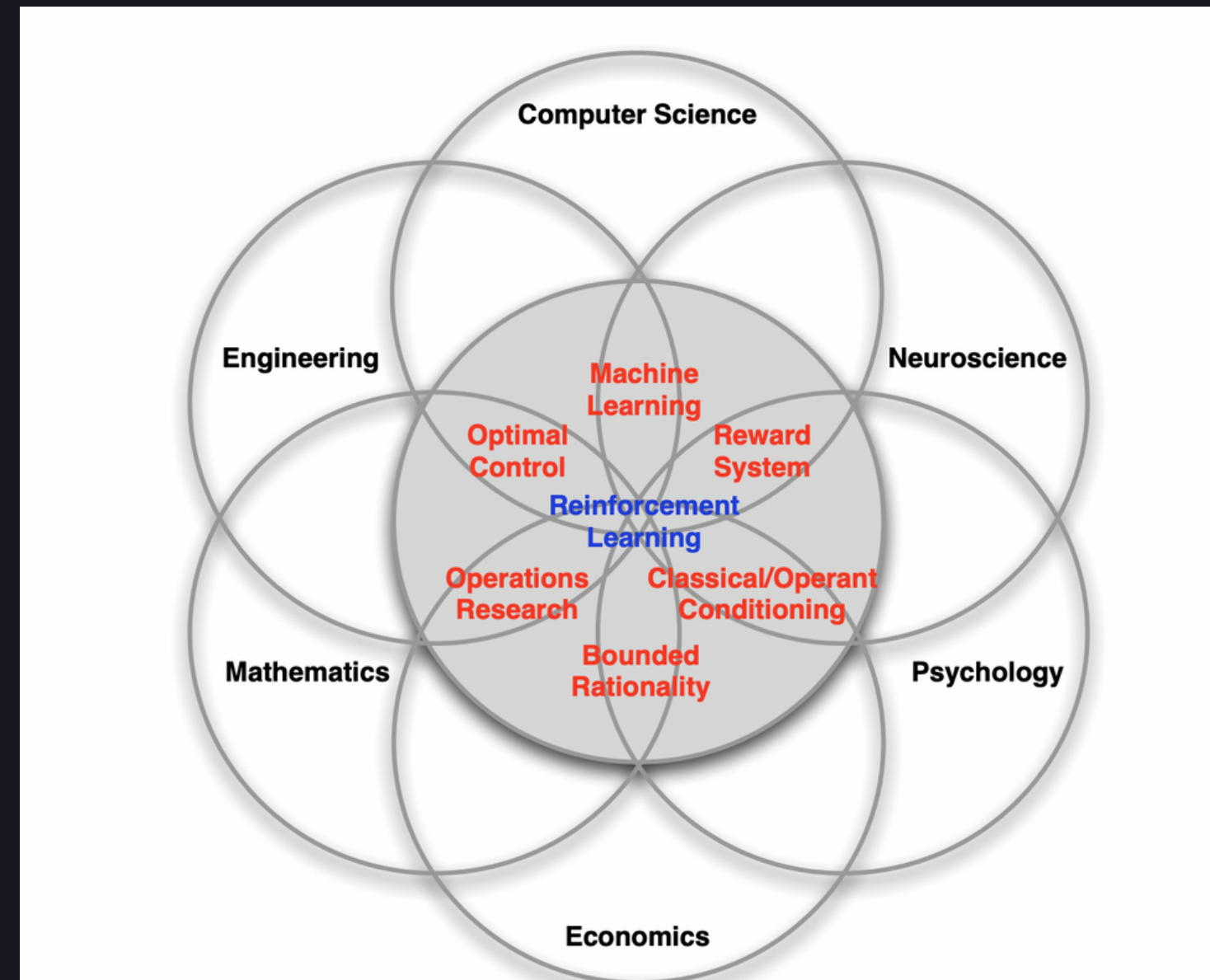background

**03** | **Policy gradient
algorithms**

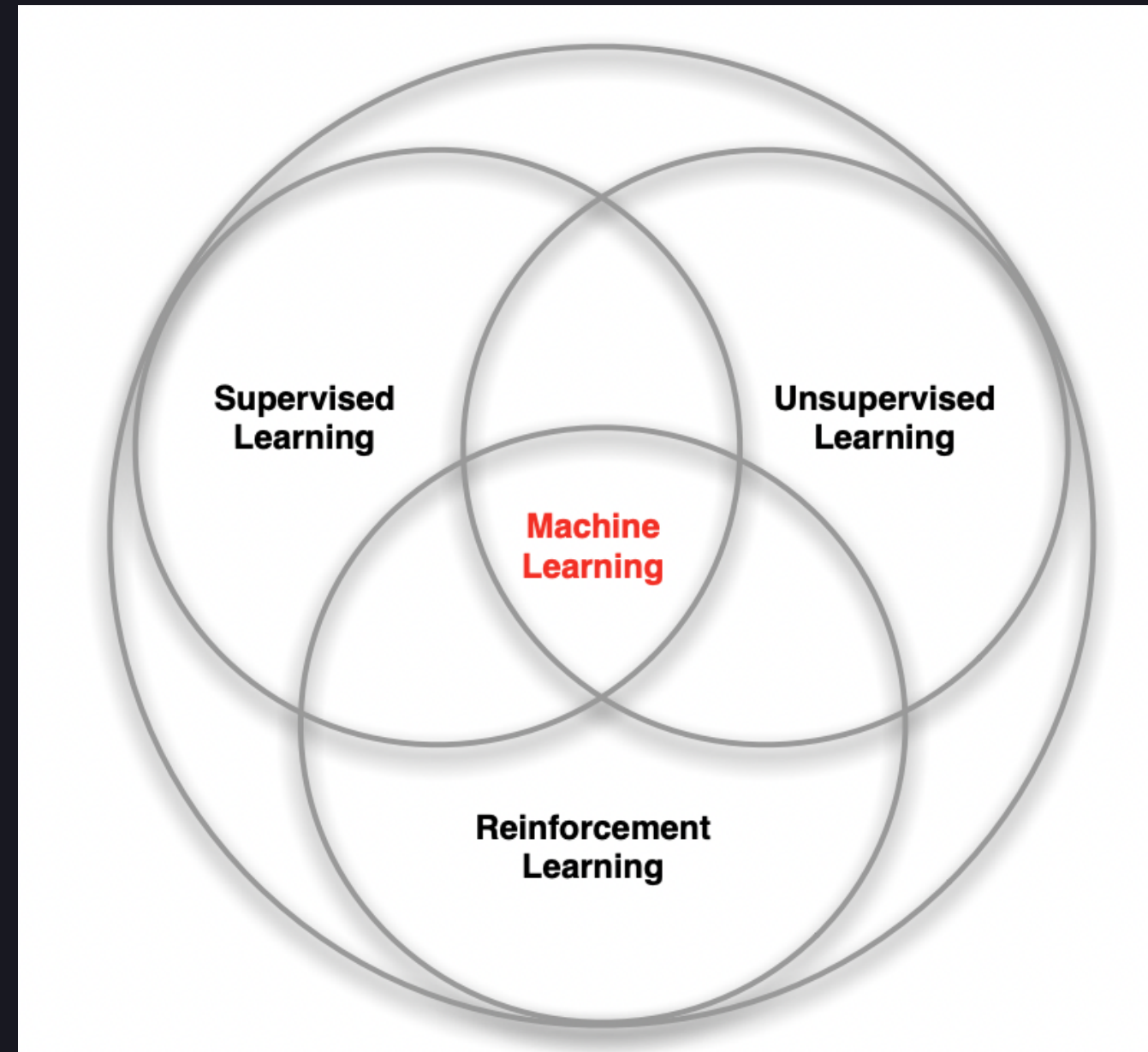**04** | **Implementation**
Using RL in practice.

# INTRO TO RL

## What is RL?
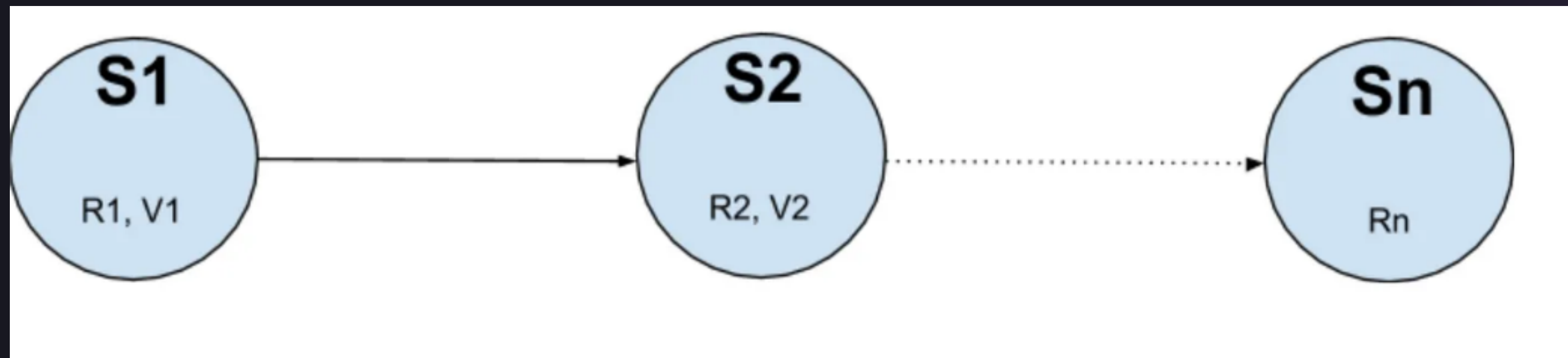
# INTRO TO RL

## What differentiates it?

# INTRO TO RL

$$V_\pi(s) = \sum_a \pi(a|s) \sum_{s'\in S} \sum_{r\in R} p(s',r \mid s,a)(r + \gamma V_\pi(s'))$$

Value function of Reinforcement Learning

# INTRO TO RL

## States and Rewards

# INTRO TO RL

## States and Rewards

$$V(s) = \sum_t \gamma^t R(S_t)$$

# INTRO TO RL

## Policy Gradient Algorithm

The reward function is defined as:

$$J(\theta) = \sum_{s \in \mathcal{S}} d^{\pi}(s) V^{\pi}(s) = \sum_{s \in \mathcal{S}} d^{\pi}(s) \sum_{a \in \mathcal{A}} \pi_{\theta}(a|s) Q^{\pi}(s, a)$$

# INTRO TO RL

Proof

$$\nabla_\theta V^\pi(s)$$

$$= \nabla_\theta \left( \sum_{a \in \mathcal{A}} \pi_\theta(a|s) Q^\pi(s,a) \right)$$

$$= \sum_{a \in \mathcal{A}} \left( \nabla_\theta \pi_\theta(a|s) Q^\pi(s,a) + \pi_\theta(a|s) \nabla_\theta Q^\pi(s,a) \right)$$

$$= \sum_{a \in \mathcal{A}} \left( \nabla_\theta \pi_\theta(a|s) Q^\pi(s,a) + \pi_\theta(a|s) \nabla_\theta \sum_{s',r} P(s',r|s,a)(r + V^\pi(s')) \right)$$

$$= \sum_{a \in \mathcal{A}} \left( \nabla_\theta \pi_\theta(a|s) Q^\pi(s,a) + \pi_\theta(a|s) \sum_{s',r} P(s',r|s,a) \nabla_\theta V^\pi(s') \right)$$

$$= \sum_{a \in \mathcal{A}} \left( \nabla_\theta \pi_\theta(a|s) Q^\pi(s,a) + \pi_\theta(a|s) \sum_{s'} P(s'|s,a) \nabla_\theta V^\pi(s') \right)$$

$$\nabla_\theta V^\pi(s) = \sum_{a \in \mathcal{A}} \left( \nabla_\theta \pi_\theta(a|s) Q^\pi(s,a) + \pi_\theta(a|s) \sum_{s'} P(s'|s,a) \nabla_\theta V^\pi(s') \right)$$

THANK YOU 😎

# RL IMPLEMENTATION



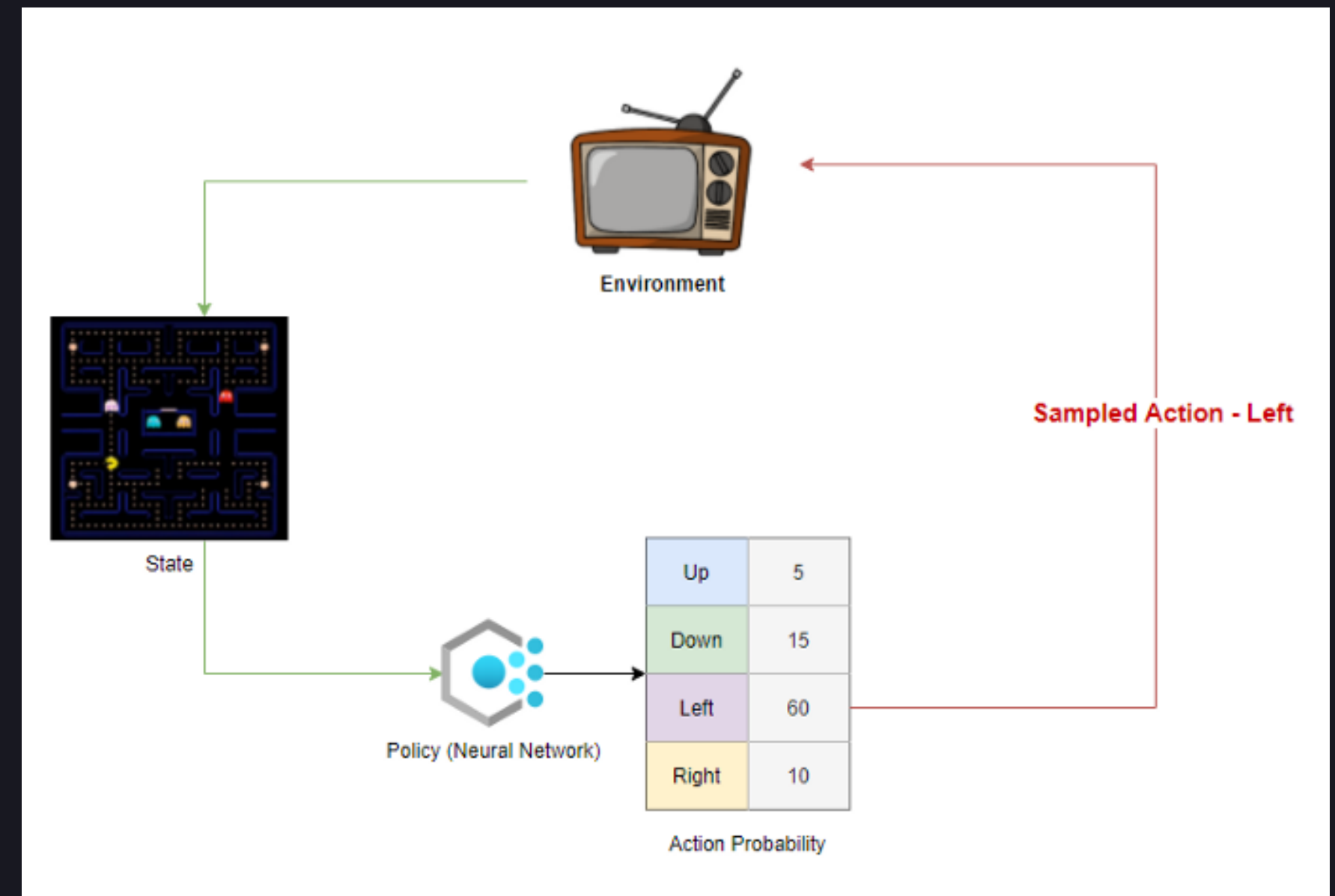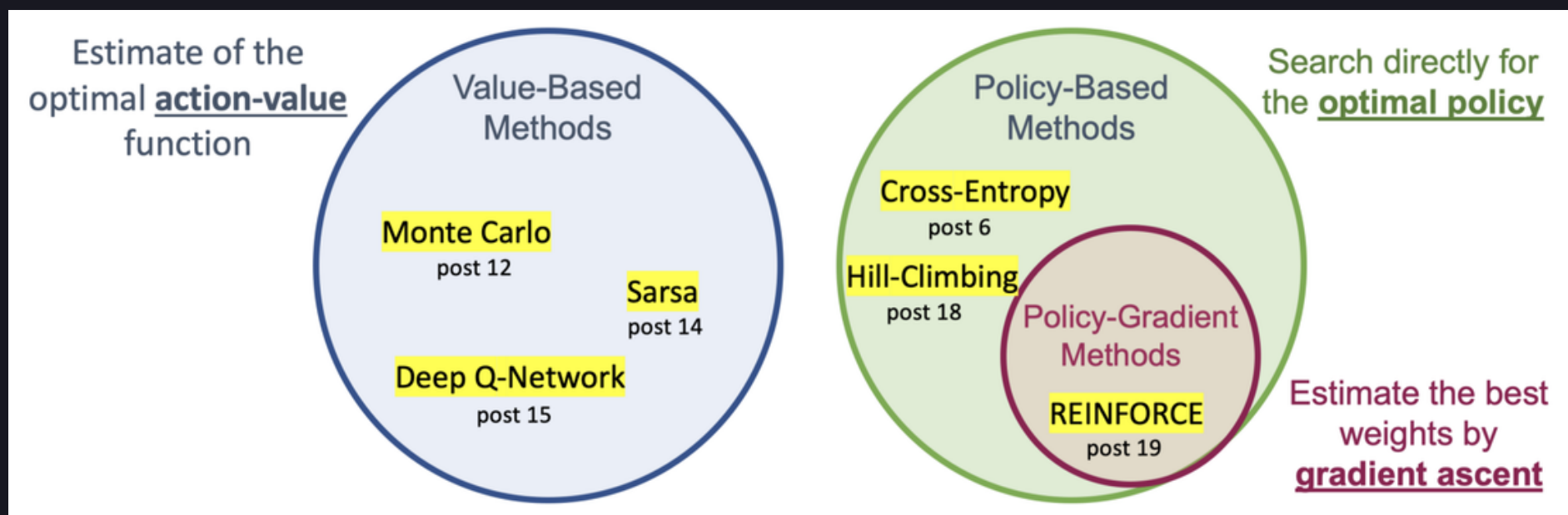*Build and train a simple agent*

# RL IMPLEMENTATION

## REINFORCE method

- REINFORCE is a policy gradient method in RL that updates policy based on expected rewards.

- It computes gradient of expected cumulative reward with respect to policy parameters.

- It is computationally efficient and used to learn stochastic/deterministic policies, but suffers from high variance and slow convergence.

# RL IMPLEMENTATION

## REINFORCE method



Estimate of the optimal **action-value** function

Value-Based Methods

Monte Carlo
post 12

Sarsa
post 14

Deep Q-Network
post 15

Policy-Based Methods

Search directly for the **optimal policy**

Cross-Entropy
post 6

Hill-Climbing
post 18

Policy-Gradient Methods

REINFORCE
post 19

Estimate the best weights by **gradient ascent**



Environment

Sampled Action - Left

State

Policy (Neural Network)

| | |
|---|---|
| Up | 5 |
| Down | 15 |
| Left | 60 |
| Right | 10 |

Action Probability

# RL IMPLEMENTATION

## Actor–critic reinforcement learning

- Actor-critic RL is a type of algorithm that combines both value-based and policy-based methods in reinforcement learning.

- It uses an actor to learn an optimal policy for decision making and a critic to estimate the value function associated with the policy.

- The actor learns to improve its policy by using the feedback from the critic, which helps to optimize the decision-making process over time.

# RL IMPLEMENTATION

## Actor–Critic

Estimates action probability

Estimates value for action

$$\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta} \left[ \nabla_\theta \log \pi_\theta(s, a) \, G_t \right] \qquad \text{REINFORCE}$$

$$= \mathbb{E}_{\pi_\theta} \left[ \nabla_\theta \log \pi_\theta(s, a) \, Q^w(s, a) \right] \qquad \text{Q Actor-Critic}$$

$$= \mathbb{E}_{\pi_\theta} \left[ \nabla_\theta \log \pi_\theta(s, a) \, A^w(s, a) \right] \qquad \text{Advantage Actor-Critic}$$

$$= \mathbb{E}_{\pi_\theta} \left[ \nabla_\theta \log \pi_\theta(s, a) \, \delta \right] \qquad \text{TD Actor-Critic}$$

SEE YOU NEXT TIME