

Reinforcement Learning

Carrot or stick?

Reinforcement Learning



Components of RL

- Agent
- Environment
- State
- Actions
- Policy
- Returns
- Exploration & Exploitation

The model that acts and learns within our environment.

The limited space that the agent exists within that defines the rules of what is possible.

The relevant information about the environment with respect to the agent's actions and results.

- State is the true reality. The agent only has access to observations.
- A state is *fully observed* if the observations represent the full state. *Partially observed* otherwise.

$$\mathbf{s}_{t+1} = f(\mathbf{s}_t, \mathbf{a}_t).$$

The agent can perform actions that alter the state.

- Discrete vs Continuous actions

The set of rules an agent follows to determine the action it should take given its observations about the environment.

Returns

The environment provides a returned value to the agent after it takes an action; reward or punishment.

We evaluate a model based on the total rewards at the end of the simulation or at time τ using,

$$R(\tau)_{finite} = \sum_{t=0}^T r_t, \quad R(\tau)_{infinite} = \sum_{t=0}^{\infty} \gamma^t r_t.$$

Exploration vs Exploitation

We want to balance between learning new things (exploration) and utilizing the knowledge we already have (exploitation).

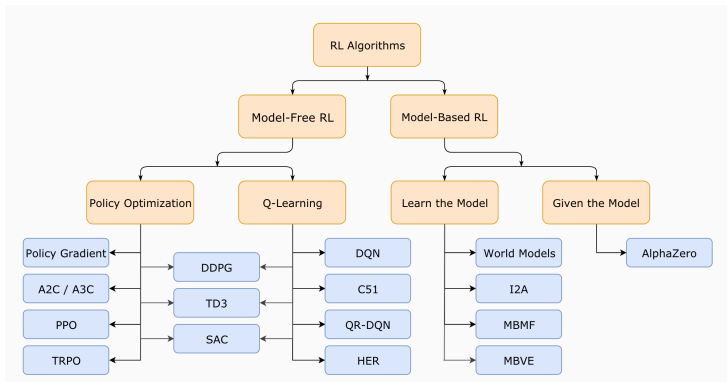
$$V^{\pi}(s) = \mathbb{E}_{\tau \sim \pi} [R(\tau) | s_0 = s],$$
$$Q^{\pi}(s, a) = \mathbb{E}_{\tau \sim \pi} [R(\tau) | s_0 = s, a_0 = a].$$

Bellman Equations

$$\begin{aligned}V^{\pi}(s_t) &= \mathbb{E} [r(s_t, a_t) + \gamma V^{\pi}(s_{t+1})], \\Q^{\pi}(s_t, a_t) &= \mathbb{E} [r(s_t, a_t) + \gamma \mathbb{E} [Q^{\pi}(s_{t+1}, a_{t+1})]].\end{aligned}$$

- Importance of reward functions - see "Curiosity Driven Learning"
- Over fitting and learned helplessness

Types of RL



Approximate action-value function $Q(s, a)$

- Create a Q table
- Select actions and update Q values
- Repeat for many iterations
- Select highest scoring action in practice

Q Table

Actions : ↑ → ↓ ←

Start				
Nothing / Blank				
Power				
Mines				
END				

ϵ -Greedy Action Selection

$$a_t = \begin{cases} \arg \max_a Q(s_t, a) & P(1 - \epsilon) \\ \text{random action} & P(\epsilon) \end{cases}$$

$$Q_{\text{new}}(s, a) = (1 - \alpha)Q(s, a) + \alpha(r_t + \gamma \max_a Q(s_{t+1}, a)),$$

where α is the learning rate and γ is our discount rate.

Using Q Table

Given a state and action, select the highest scoring action.

Questions

These slides are designed for educational purposes, specifically the CSCI-470 Introduction to Machine Learning course at the Colorado School of Mines as part of the Department of Computer Science.

Some content in these slides are obtained from external sources and may be copyright sensitive. Copyright and all rights therein are retained by the respective authors or by other copyright holders. Distributing or reposting the whole or part of these slides not for academic use is HIGHLY prohibited, unless explicit permission from all copyright holders is granted.