

# Reinforcement Learning

Lecture - 01

## Artificial Intelligence :-

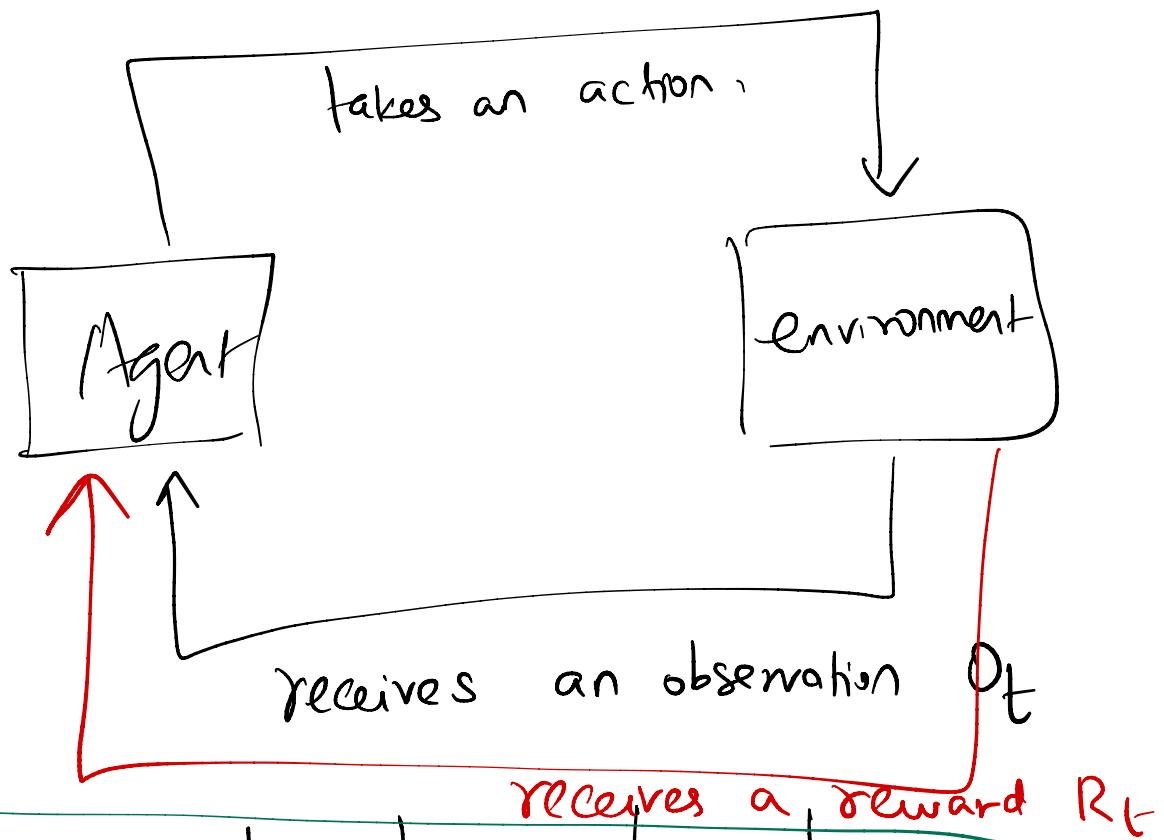
Goal of AI : → to understand the general principles behind human intelligence.

→ to simulate them in machines.

ML : takes the "learning" approach to solve AI.

→ Active learning. (RL) *Learning through interaction.*

→ Passive learning. (SL) *Learning from data.*



Objective: to learn how to map situations to action — so as to maximize a numerical reward signal.

Example: Trash Collecting robot.

+1 — whenever it collects a trash,

0 — rest of the time

-100 — if it runs out of  
battery.

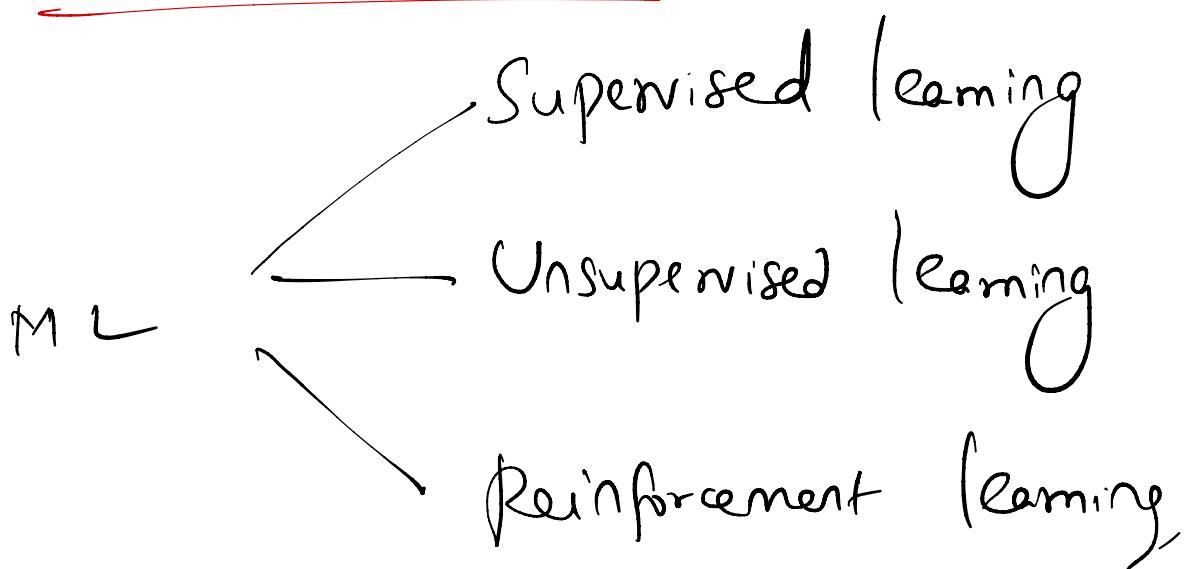
Reward hypothesis:

That all of what we mean by  
goals and purposes can be well  
thought of as max. of the  
expected value of the cumulative  
sum of a received scalar  
signal (reward)

## Reward-is-Enough Hypothesis (Silver et al. 2021)

Intelligence, and its associated abilities, can be understood as subserving the max. of reward by an agent acting in its environment.

Defour to ML :-



Supervised learning:

$$D_{Tr} = \{x^{(i)}, y^{(i)}\}_{i=1}^N$$

$$f: x \rightarrow y$$

new  $x \rightarrow ?$

$x$ : review

$y$ : + Positive review

- Negative rev

"Prediction"

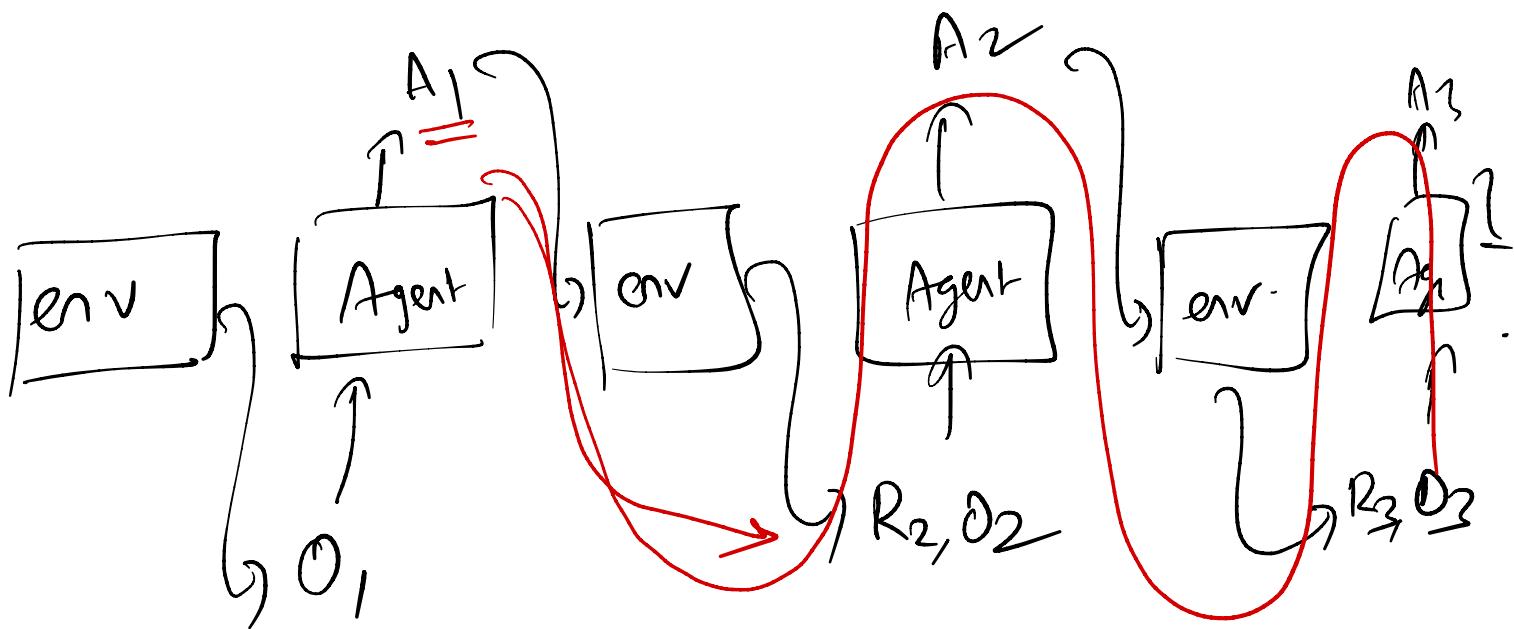
## Unsupervised learning:-

$$\mathcal{D} = \{x^{(i)}\}_{i=1}^N$$

goal: find some structure in the data.

ex: clustering

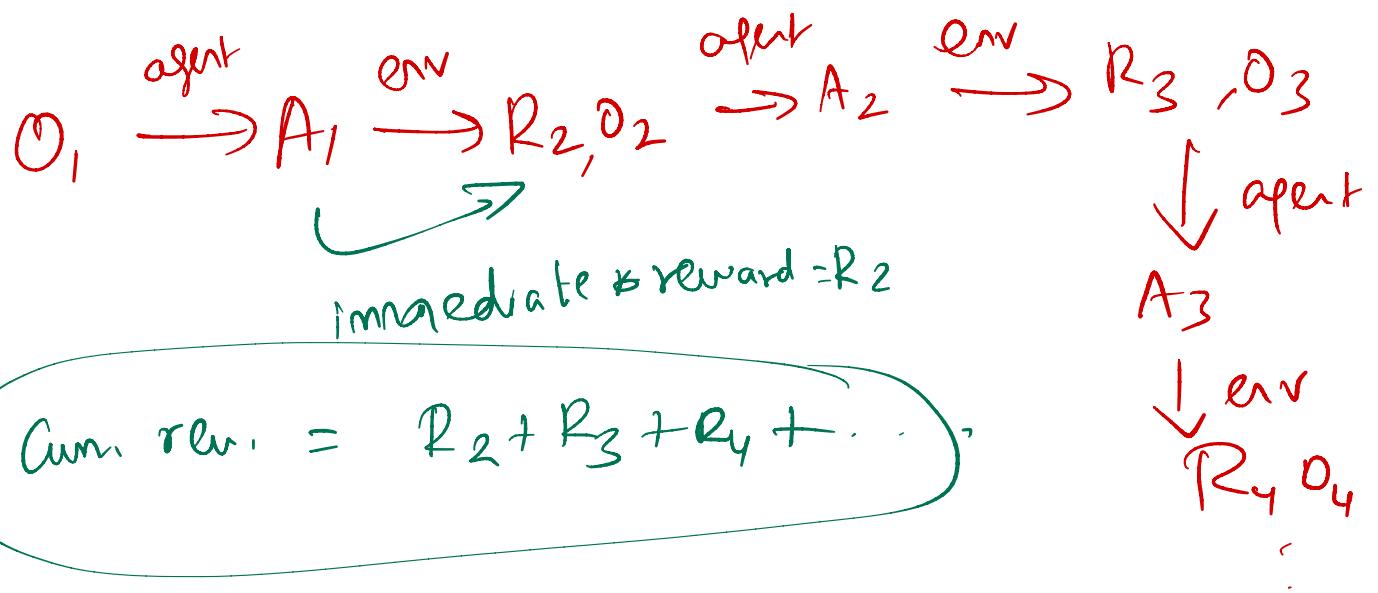
# Reinforcement Learning:



- Interaction
- notion of "time"

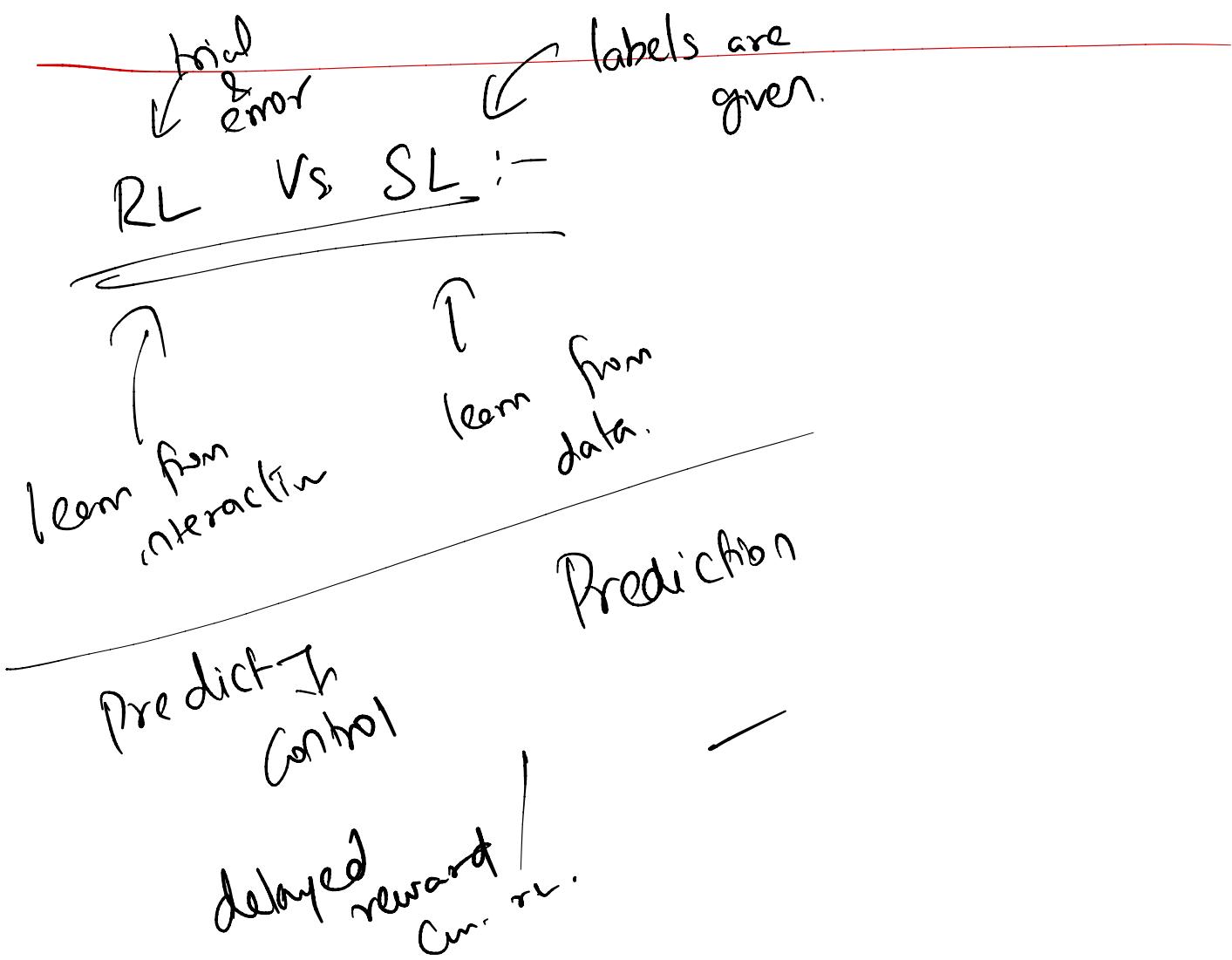
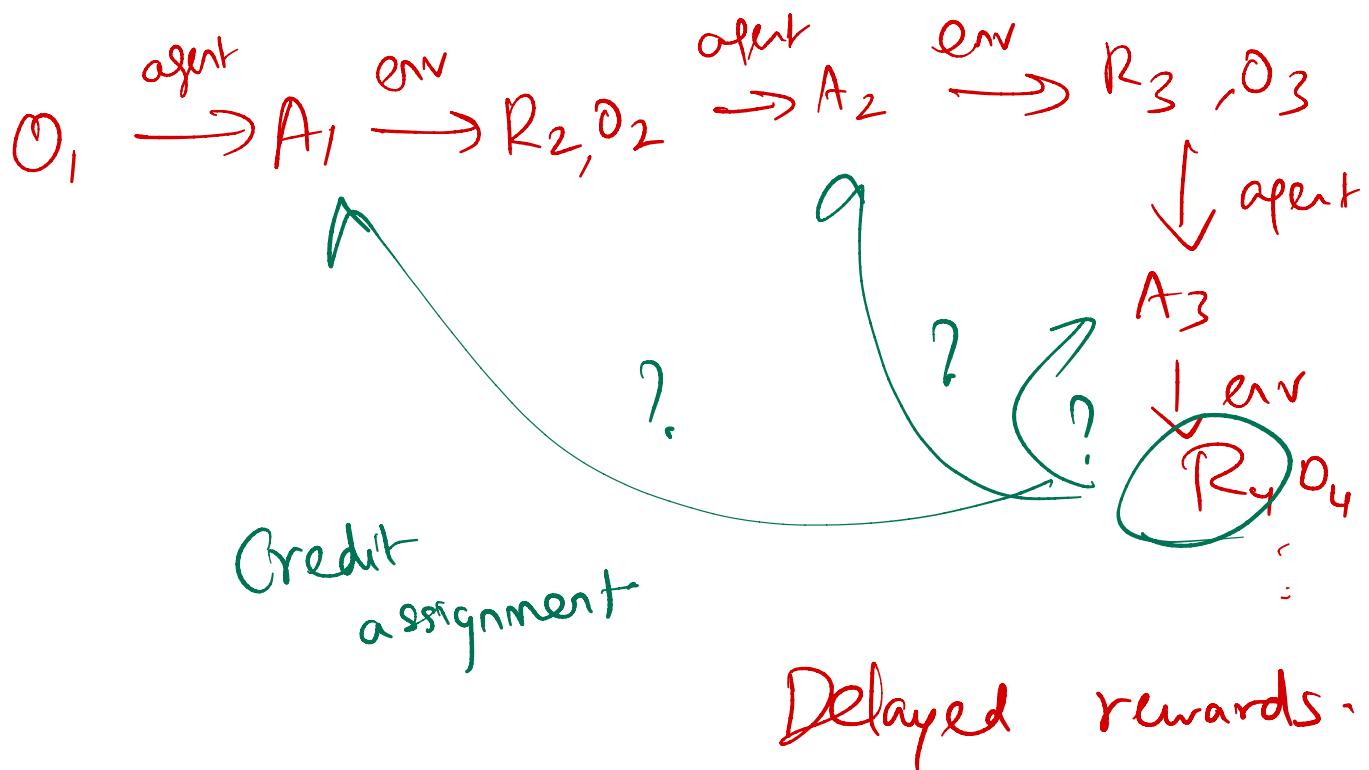
## Implications of Closed Loop:-

Forward view:-



goal: to max.  
cumulative reward

Backward view:



RL vs. VS L :-

Exploration vs. Exploitation dilemma :-

Exploit  $\rightarrow$  what it has already exp.  
in order to obtain reward

Explore  $\rightarrow$  in order to make better  
action selection in the  
future.

## RL - Summary:-

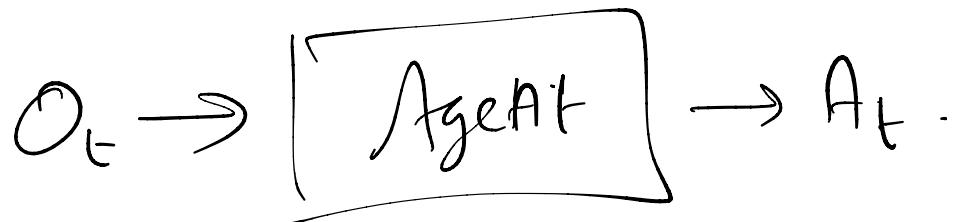
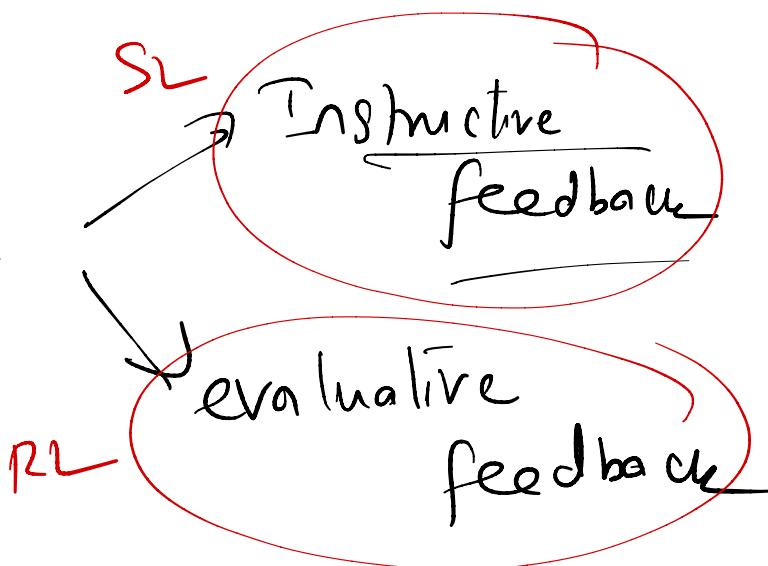
- ① Agent view of AI.
- ② Sequential decision making process.
- ③ Delayed rewards & credit assignment
- ④ Exploration / Exploitation dilemma.

## Immediate RL

Key char. of an RL problem:-

- ① Learning to act in many situations.
- ② Delayed rewards / credit assignment
- ③ Exploration / Exploitation dilemma.

types of feedback



$(A_1, A_2, A_3, A_4)$   
 $a^* = A_2$

Optimal action :  $a^*$

Immediate RL: / k-armed bandit problem

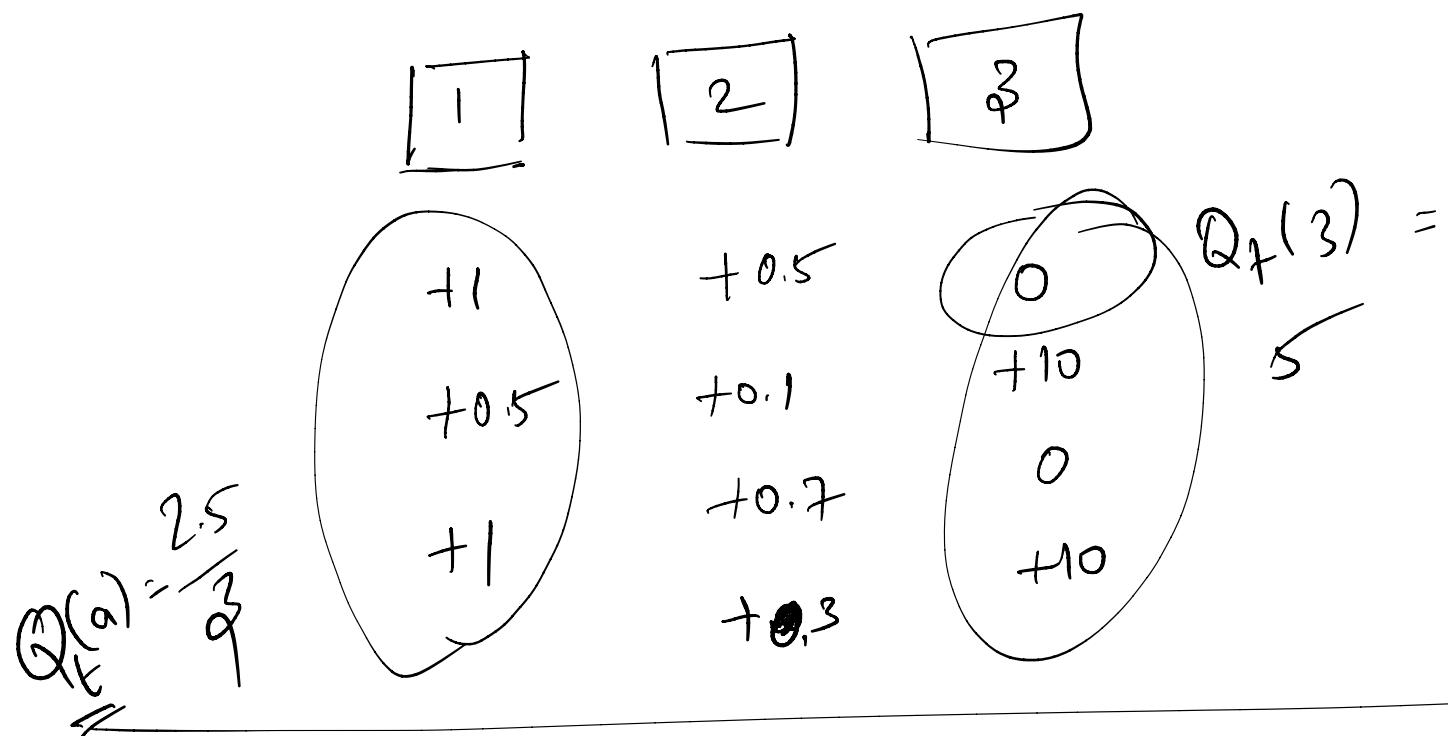


options/actions.

reward  $\rightarrow$  chosen from stationary  
Prob. distrib-  
that depends on your action.

Objective: to maximize the expected total reward over some time period.

- Note:
1. Agent sees the same state (some set of actions) all the time.
  2. Rewards are immediate



Value of action  $a = q_{t+}(a)$

$$= E[R_t \mid A_t = a]$$

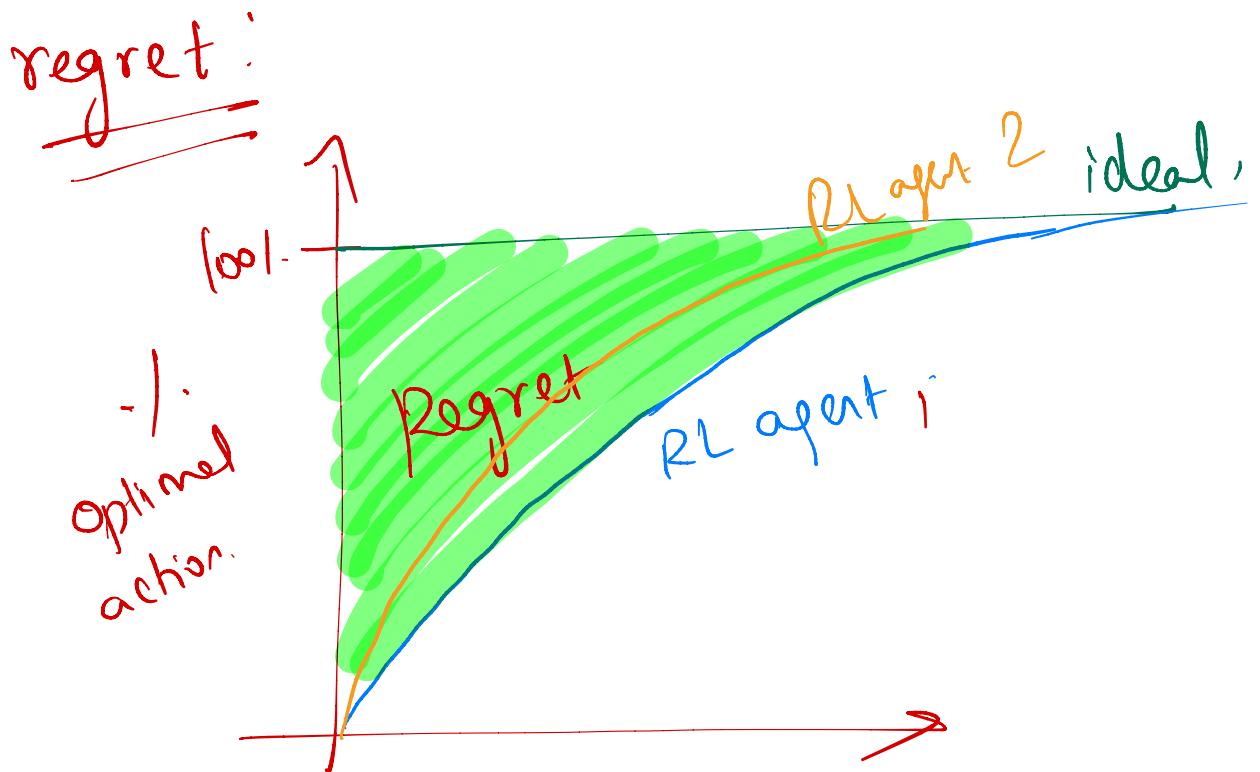
$q_{t+}(a)$   
actual.

$Q_t(a)$   
estimate

$$a^* = \operatorname{argmax} q_{t+}(a)$$

$a \in \{a_1, \dots, a_k\}$

$$A_t = \operatorname{argmax}_{a \in \{a_1, \dots, a_k\}} Q_t(a)$$



$$\text{Regret} = K q^*(a^*) - \sum_{t=1}^K R_t$$

## Action - Value methods:-

$Q_t(a)$  = Sum of rewards when 'a' taken  
prior to 't'

# of times 'a' taken prior to 't'.

$$Q_t(a) = \frac{\sum_{i=1}^{t-1} R_i \mathbb{1}_{A_i=a}}{\sum_{i=1}^{t-1} \mathbb{1}_{A_i=a}}$$

$$\mathbb{1}_{A_i=a} = \begin{cases} 1 & \text{if action } 'a' \text{ was taken} \\ & \text{in time step } 'i' \\ 0 & \text{otherwise} \end{cases}$$

As  $t \rightarrow \infty$   $Q_t(a) \rightarrow q_{\pi^*}(a)$

Exploit :  $A_t = \underset{a}{\operatorname{arg\max}} Q_t(a)$

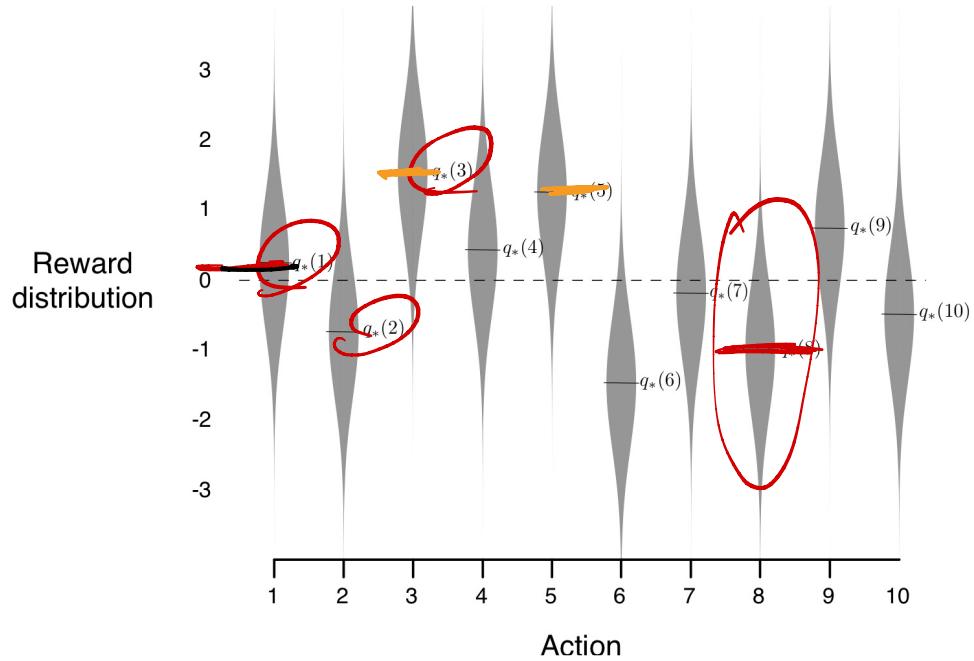
"greedy" action.

Prob.  $\epsilon$        $0.1, 0.01, 0.05$

$$A_t = \begin{cases} \underset{a}{\operatorname{argmax}} Q_t(a) & \epsilon \\ \text{random action} & 1 - \epsilon \end{cases}$$

$\epsilon$ -greedy





$$\text{mean } \alpha_x(i) = \mathcal{N}(0, 1)$$

