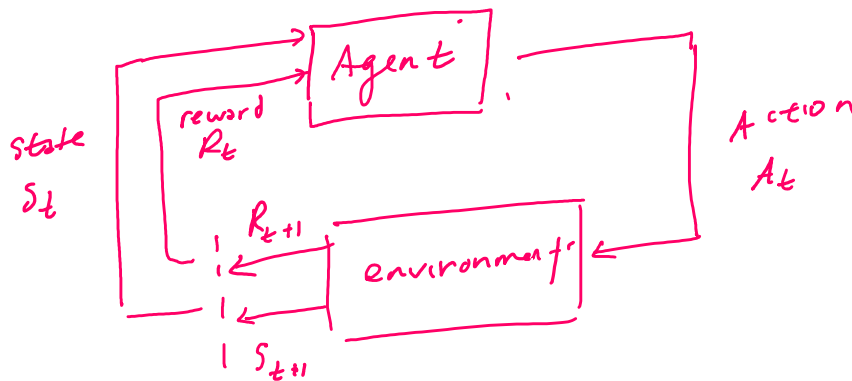# Summary of The RL Framework: The Problem

Saturday, November 28, 2020    1:28 PM



## Setting

- RL Framework: Agent learning with an environment
- Each step: → Agent recieves environment's state
  → Agent performs action based on state
- Next step: → Reward recieved from last step
  → New environment state recieved
  → Agent performs next action

## Episodic vs. Continuing Tasks

- task: instance of RL Framework
- Continuing Tasks: Continue forever
- Episodic Tasks: well-defined starting & ending
  b) under when agent reaches terminal state

## Reward Hypothesis : All goals are aimed to maximize expected culimative reward.

Culmative Reward :  $G_t \doteq R_{t+1} + R_{t+2} \ldots$

Discounted Return:

$$G_t := \gamma^k \sum_{k=0}^{\infty} R_{k+1} \quad , \text{where} \quad \gamma \in [0,1]$$

$\hookrightarrow \gamma \Rightarrow 0$    immediate reward

$\gamma \Rightarrow 1$ , long-term

if $\gamma = 1$ , no discount

## MDP (Merkov Decision Process):

$\rightarrow$ state spaces all non-terminal states ( S )

$\hookrightarrow S^+$ includes terminal states

$\rightarrow$ action space : set of possible actions ( A )

$\rightarrow$ one-step Dynamic:

$$P(s', r | s, a) = P(S_{t+1} = s', R_{t+1} = r | S_t = s, A_t = a)$$

$\rightarrow$ Finite MDP:

- finite S
- finite A
- set of Rewards
- one-step dynamics of environment
- discount rate $\gamma \in [0,1]$