

Lecture 2

RL: Learn by interacting with the environment

- more close to the nature of learning

(Compared with supervised learning)

- no explicit teacher

- Discover which action yield the most reward by trying it

$$S_t \xrightarrow{a_t} S_{t+1}, r_{t+1}$$

↑ ↑
stochastic stochastic & delayed.

Formulation:

\mathcal{S} ,

\mathcal{A} ,

\mathcal{P} , $P(S_{t+1} | S_t, a_t) \leftarrow$ distribution of the next state given the current state & action

$r: \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$

γ : discount factor $\gamma \in (0, 1]$

Goal: $\max \mathbb{E} \left[\sum_{i=0}^{\infty} \gamma^i r_i \mid S_0 = s \right]$

Difficulty: - γ close to 1, take more future reward into account

- The optimal action may not be the one who has the highest immediate reward

- After I apply an action, do not know what the next state is.

- the reward could be a random variable

RL: $\mathcal{S}, A, \underset{\text{unknown}}{P}, r, \gamma$

MDP: $\mathcal{S}, A, \underset{\text{known}}{P}, r, \gamma$

MAB: A, r

- UCB \rightarrow Thm 1: upper bd for regret
- Bayes-optimal

LP: $\mathcal{S}, A, \underset{\text{finite}}{\overbrace{P}^{\text{finite}}}, r, \gamma$

primal, Dual, primal-dual
 \downarrow Bellman \downarrow policy ∇

\rightarrow Thm 2
 \uparrow
Equivalence.

RL: TD \rightarrow Thm 3: convergence

- GTD, policy ∇ , AG,
- HJB: continuous-time control with known dynamics.
- HJB \Leftrightarrow optimal control \rightarrow Thm 4.