

Unit 1

- What is RL

- big picture:

- * *Agent(AI)* will *Learn* from *Environment* by *Interacting* with Environment and receive *Rewards* (neg/pos) as *Feedback* for performing *Actions*

- * while not done:

- Environment -> observation -> Agent

- Agent -> action -> Environment

- Environment -> reward, new observation, done -> Agent

- example:

- boy learning to play video game

- boy is Agent,

- screen is observation

- buttons is set of possible actions

- game is environment

- rewards is points

- boy observes screen (sees avatar, coin, squid)

- boy presses right button

- game reacts by updating screen -> avatar touches coin, rewards point

- boy learns getting coin rewards +1

- boy presses button again and tough

- game reacts by updating screen -> avatar touches squid and dies, reward

- boy learns touching squid rewards dead, and ends game

- formal definition:

- * RL is **FRAMEWORK** for solving **CONTROL TASKS** (aka decision problems)

- * by building **AGENTS**

- * AGENTS that **LEARN** from the **ENVIRONMENT**

- LEARN by **INTERACTING** with ENVIRONMENT

- LEARN thru **TRIAL** and **ERROR**

- LEARN by receiving **REWARDS**(NEGATIVE/POSITIVE) as **FEED-BACK**
- RL Framework
 - RL Process
 - * Environment \rightarrow (state S_t , reward R_t) \rightarrow Agent \rightarrow action $A_t \rightarrow$ Environment \rightarrow (state S_{t+1} , reward R_{t+1}) \rightarrow Agent (loop)
 - * IOW: $S_0 \rightarrow A_0 \rightarrow R_1, S_1 \rightarrow A_1 \rightarrow R_2, S_2 \dots R_n, S_n \rightarrow A_n \rightarrow R_{n+1}, S_{n+1} \rightarrow \dots$
 - * Agent's GOAL: Maximize its **CUMULATIVE REWARD** aka **EXPECTED RETURN**
 - * WHY IS THIS(Maximization of EXPECTED RETURN) the Agent's goal? b/c RL is based the REWARD HYPOTHESIS
 - the central idea of RL: The reward hypothesis
 - * ALL GOALS can be described as the MAXIMIZATION of EXPECTED RETURN
 - Example?
 - * TO HAVE BEST BEHAVIOR, MAXIMIZE the EXPECTED CUMULATIVE REWARD
 - are EXPECTED CUMULATIVE REWARD same as EXPECTED RETURN? \Rightarrow YES!
 - ANOTHER NAME for RL Process: MARKOV DECISION PROCESS (MDP)
 - Markov property AGENT only needs CURRENT STATE to decide what ACTION to TAKE - CONTRAST with NEEDING HISTORY of ALL STATES and ACTIONS they took before
 - Observations/States Space Observations/States - information our AGENT gets from the ENVIRONMENT. - Example: videogame, \rightarrow Observation/State is a Screenshot (AKA Frame)
 - Action Space
 - Rewards and the discounting
 - Task Types
- Exploration/Exploitation Tradeoff
- Solving RL Problems: 2 main approaches
 - The Policy PI
 - Policy based Methods
 - Value based Methods
- “Deep” in Deep RL

$$z(x) = \sum_{i=1}^n w_i x_i + b = w \cdot x + b$$