# Reinforcement Learning

→ Science of decision making
→ other fields have other names for it

RL ≠ from other methods:

- → no supervisor, only rewards
- → feedback is delayed
- → time really matters (sequential, no i.d. data)
- → agent has its own actions → influencing the environment
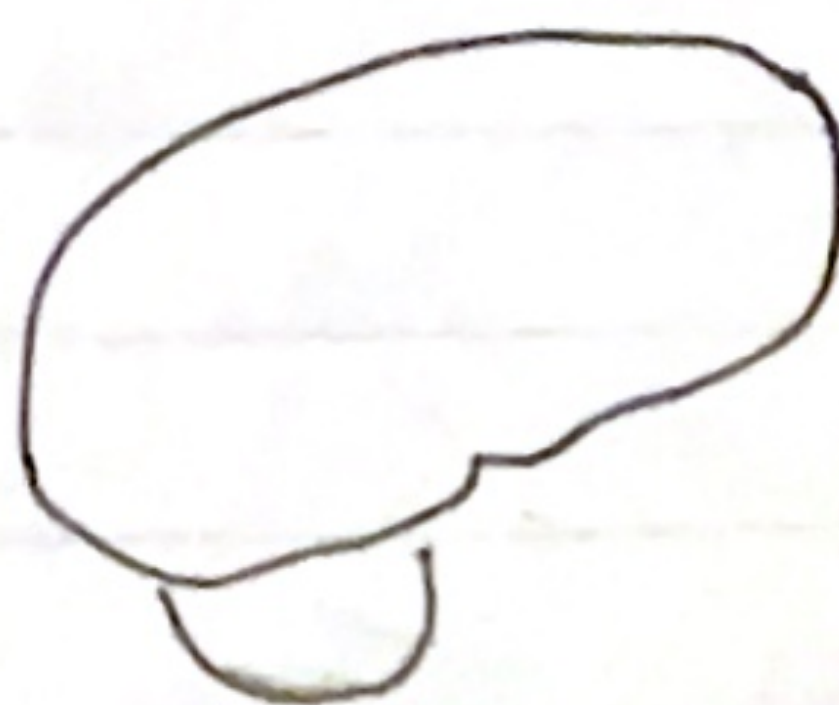  - → influencing the data it receives —

Reward hypothesis → they look for the maximization of the cumulative reward

Goal — select actions to maximize the future reward.
Actions might have long term consequences



Observation
$\xrightarrow{\quad}$
$o_t$

→ Action
$A_t$

↑
Reward → Scalar
$R_t$

history sequence of observations = $H_t = A_1, O_1, R_1, \ldots A_t, O_t, R_t$

State → concise summary of the history

↓

environment state ⌐ agent doesn't see them but
the algorithm cannot depend on it
good idea to manipulate it (doesn't matter) to make it easier

Agent state → $s^a$ in our algorithm — used to pick next action
↳ our decision — what to use / what to throw away

Information State (Markov State) contains all useful information from history

A State $S_t$ is Markov if and only if:

$$\mathbb{P}[S_{t+1} \mid S_t] = \mathbb{P}[S_{t+1} \mid S_1, \ldots, S_t]$$

probability of next ___ is the same if you show
state, conditioned ___ all past info
the state you're in

⌣

You can throw away all your past
info if you have your present info

"the future is independent of the past given the present"

$$H_{1:t} \to S_t \to H_{t + 1:\infty}$$

↓

State is sufficient statistic of the future

→ What we believe will
representation of state

Fully observable environment
state

$O_t = S$

Agent state = environment

Markov

Partial observability → indirect
(eg robot with ca

|

Partially Obs

Agent must construct its
• beliefs of envir
• complete history
• recurrent neural n

And RL might include:
policy → agent's behaviour
value function → how
model → agent's

→ What we believe will happen next, depends on our representation of state (eg lever and cheese)

Fully observable environment = agent directly observes environment state

$$O_t = S_t^a = S_t^e \quad (?)$$

Agent state = environment state = information state

↓

Markov Decision Process
(MDP)

Partial observability → indirectly observes environment
(eg robot with camera)

|

Partially Observable MDP (POMDP)

Agent must construct its own state representation
- beliefs of environmental state
- complete history
- recurrent neural network

And RL might include:
policy → agent's behaviour function
value function → how good is each state and/or function
model → agent's representation of environment

Policy —

mag from state to action

⌈ Deterministic → $a = \pi(s)$
⌊ Stochastic $\pi(a \mid s) = \mathbb{P}[A = a \mid S = s]$

→ taken particular action given
certain state

Categorizing RL agents

Value based
Used to evaluate goodness / badness of State

$$V_\pi(s) = \mathbb{E}_\pi[R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \cdots \mid S_t = s]$$

Model based

⌈ predictions — P predicts the next state (dynamics)
⌊ rewards $\qquad P_{ss'}^a = \mathbb{P}[S' = s' \mid S = s, A = a]$

$R^a$ predicts the next (immediate) reward

$$R_s^a = \mathbb{E}[R \mid S = s, A = a]$$

Actor critic — both ⌈ model
                        ⌊ value

(2)

Model free → policy and / or free value function

Model based

va

# Key Subproblems

→ Learning vs Planning → environment is unknown and agent learns through interaction. In planning, the environment model is fully known, and the agent performs internal look-ahead search

| RL often combines both

1° learning a model
2° planning with it

→ Exploration vs Exploitation
  → Universal tradeoff unique to RL
  → exploitation: use current knowledge to maximize known reward
  → exploration: sacrifices immediate reward to discover better long-term option
  (eg trying new restaurant vs favourite restaurant)

→ Prediction vs Control
  Prediction answers "how well will I do following my current policy?"
  Control asks "what is the optimal policy?"

  → Solving prediction problem is typically necessary for solving control