

Intelligence := to be able to learn to make decisions to achieve goals

What is RL?

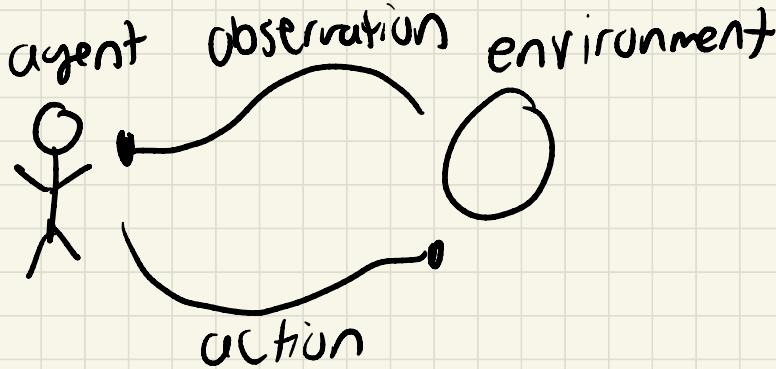
- People/animals learn by interacting with env.

↳ active

↳ sequential

- We are goal-directed

- We can learn without examples of optimal behaviour. → we optimise reward signal



Goal : optimise sum of rewards, through repeated interaction

RL is based on the reward hypothesis!

Any goal can be formalized as the outcome of maximizing a cumulative reward.

Examples of RL

- fly helicopter
- managed investment portfolio
- control power station
- make a robot walk
- play video or board games.

If the goal is to learn via interaction, it is RL.

2 distinct reasons to learn:

RL can do both.

1) Find solutions

2) Adapt online, deal with unforeseen circumstances.



RL is

- Science and framework of learning to make decisions from interaction.

↳ need to think about

- time
- long term consequences of action
- effectively gathering experience
- predicting the future
- dealing with uncertainty

- Huge potential scope

- A formalisation of the AI problem

formalising the RL Problem

- At each step t , the agent:
 - Receives Observation O_t (and reward R_t)
 - Executes action A_t
 - The environment:
 - Receives action A_t
 - Emits observation O_{t+1} (reward R_{t+1})
- $\left. \begin{matrix} \\ \\ \end{matrix} \right\}$ could be continuous time.

Rewards → Scalar feedback signal

- Indicates how well agent is doing at step t
- agent goal is to maximize cumulative reward
↳ return. (future reward)
 $G_t = R_{t+1} + R_{t+2} + R_{t+3} + \dots$

Values

- expected cumulative reward, from a state s ,
the value \rightarrow goal = maximize this.

$$\begin{aligned} v(s) &= E[G_t | S_t = s] \\ &= E[R_{t+1} + R_{t+2} + \dots | S_t = s] \end{aligned}$$

- the value depends on action agent takes
- Rewards and values define utility of states and actions
- Returns can be defined recursively

$$G_t = R_{t+1} + G_{t+1}$$

$$v(s) = E[R_{t+1} + v(S_{t+1}) | S_t = s]$$

- Goal : Select actions to maximise value.

Mapping from states to actions := policy

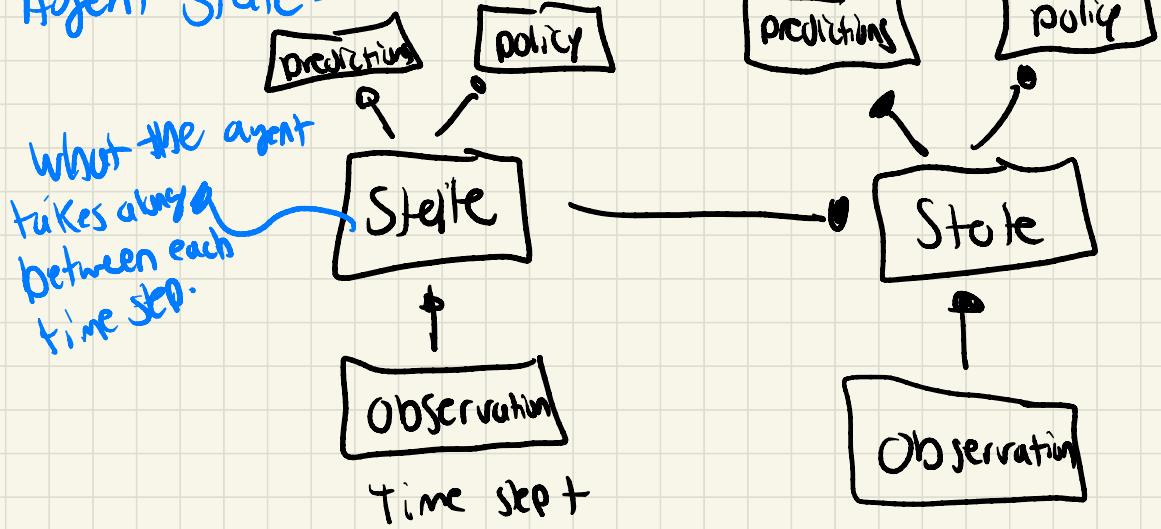
$$\begin{aligned} q(s,a) &= E[G_t | S_t = s, A_t = a] \\ &= E[R_{t+1} + R_{t+2} + \dots | S_t = s, A_t = a] \end{aligned}$$

Core concepts!

- Environment
- Reward
- Agent
 - agent state
 - policy
 - value function estimate?
 - model?

Inside the agent : the Agent State

Agent state:



The history := full sequence of observations/actions/rewards

$$H_t = O_0, A_0, R_1, O_1, \dots, O_{t-1}, A_{t-1}, R_t, O_t$$

↳ used to construct the agent state S_t

Pull observability := Agent sees full environment state (we could use environment state as state)

$S_t = O_t$ = environment state.

Markov decision processes (MDPs)

Decision process is Markov if:

$$P(r_t | s_t, a_t) = P(r_t | s_t, H_t, A_t)$$

(state contains everything we need
to know from history.)

→ once state is known, the history may
be thrown away.

Typically, s_t is some compression of H_t

Partial observability: observations are not Markovian
using observation as state would not be
Markovian, this is called partially observable
Markov decision process (POMDP)

Agent State \rightsquigarrow actions depend on state
↳ function of history

Generally: $S_{t+1} = u(S_t, A_t, R_{t+1}, O_{t+1})$
↳ state update function

To deal with partial observability, agent can construct suitable state representations

examples:

- Last observation: $S_t = o_t \rightarrow$ not enough

- Complete history: $S_t = H_t \rightarrow$ 2 large

- general update: $S_t = u(S_{t-1}, A_{t-1}, R_t, O_t)$

↳ how to pick?

Fully Markovian agent state not always feasible
more important: state allow good policies and value prediction

Inside the agent) the policy

Policy: defines agent's behavior

↳ map from agent state to action

Deterministic policy: $A = \pi(S)$

Stochastic policy: $\pi(A|S) \geq p(A|S)$

Inside the agent! Value estimates

the actual value function is the expected return

$$V_{\pi}(s) = E[G_t \mid S_t=s, \pi]$$

$$= E[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \mid S_t=s, \pi]$$

Discount factor

$$\gamma \in [0, 1]$$

- value depends on policy
- Used to evaluate the desirability of a state
 - ↳ Trades off importance of immediate vs long-term rewards

Optimal Value:

$$V_{*}(s) = \max_q E[R_{t+1} + \gamma V_{*}(s_{t+1}) \mid S_t=s, A_t=a]$$

↳ does not depend on policy

↳ useful to create algorithms.

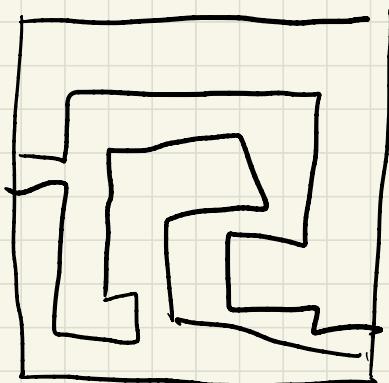
Inside the agent! Model

↳ could be stochastic.
↳ predicts what env. will do

$$p(s_t, a_t, s') \approx p(s_{t+1} = s' | s_t = s, a_t = a)$$

↳ probability of observing a
next state given a state and
action.

Example



Rewards: -1 per time step

Actions: N, E, S, W

States: Agent location

Agent Categories

1) Value Based

- No Policy
- Value Function



2) Policy Based

- Policy
- No value function

3) Actorcritic

- Policy
- Value function

- Model Free := No model

- Model Based := - Optionally policy and/or value function
- model

All components are functions

- Policies $\pi: S \rightarrow A$

- Value functions $v: S \rightarrow \mathbb{R}$

- models: $m: S \rightarrow S$ and/or $r: S \rightarrow \mathbb{R}$

- State update: $u: S \times O \rightarrow S$

(Can use neural networks and use deep learning to learn)

