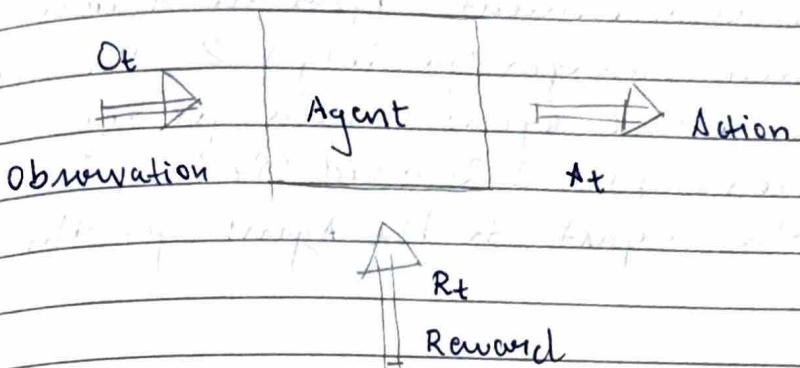
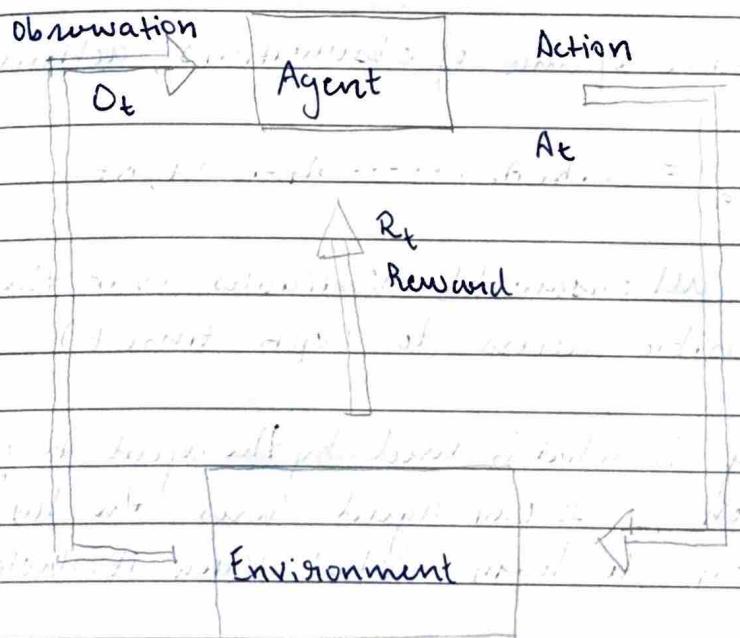


Reinforcement Learning



An agent performs an action based on the inputs it receives.

Namely O_t which is the environment or a piece of it, or a scalar feed back R_t .
Our goal is to find the algorithm for our Agent based on which it can take actions (based on information it gets) at each time step



Our Agent can only influence the environment via actions, beyond that the agent has no control over the environment. After an action, the environment changes (possibly) & the feedback is given to our agent in the form of Observation (O_t) & Reward (R_t). This again forms the input to the agent for the next Action.

So at each step t the Agent :

→ receives O_t & R_t , & based on this it takes an action A_t

And the Environment : And last action is set to

Replies the action A_t & it emits Observation O_t & Reward R_t ~~and last action is set to~~

c) History is the sequence of observations, actions, rewards

$$H_t = O_0, R_0, A_0, \dots, A_{t-1}, O_t, R_t$$

These are all observable variables that the agent has possible access to (upto time t)

So history is what is used by the agent to predict the next action, so our agent uses the history in some form to learn what action it should take in the future

Since histories tend to blowup & are large, we use States, a state is simply a (concise) summary of the history & this is what our agent uses instead of the entire history.

Formally a state at time step t is a function of history

$$s_t = f(H_t)$$

D) Types of States

Environment State s_t^e is the environment's private representation & the environment chooses O_t & R_t from s_t^e to give to the agent. This isn't visible to the agent, even if it contains irrelevant information.

Agent State s_t^a is the agent's internal representation, it summarizes everything its seen so far & uses that to pick the next action, so we simply can choose what to process, what to retain from previous time steps in our agent to help it make sure to take an action.

Also it can be any function of the history

$$s_t^a = f(H_t)$$

& we can build out this function, depending on what we want our agent to be

E) Another Mathematical definition of State:

An information or Markov State is one which contains all useful information from the history.

Also a Markov Property for a State s_t is

$$P[s_{t+1} | s_t] = P[s_{t+1} | s_1, \dots, s_t]$$

The property says that the probability of the next state (s_{t+1}) conditioned on the previous state s_t is the same for the next state even if it is conditioned on all the previous states (s_1, \dots, s_t).

So s_t contains enough information that discarding all the past we can still make decisions about the future simply by giving the present state s_t .

"Future is independent of the past given the present"

$$H_{1:t} \rightarrow s_t \rightarrow H_{t+\infty}$$

Given this State s_t , we can throw away the history & hopefully this state is a sufficient statistic of the future. s_t fully characterizes the distribution over future actions & rewards.

To understand a markov state, s_t to be able to discard previous states & still be able to take decisions about the future, so such a state s_t must comprise enough information to be able to do this.

This an environment state s_t^e is markov, since the full access to an environment's information is enough to characterize future states regardless of the history

Trivially the history is also markov, since the History itself is the culmination of all previous States (Observations), Rewards & Actions.

F) Fully Observable Environments

The Agent directly observes the (entire) environment state s_t^e

so the quantities

$$o_t = s_t^a = s_t^e \text{ all happen to be the same}$$

Formally this is called the Markov Decision Process (MDP)

G) Partially Observable Environments

The agent does not get to see the environment directly, entirely but rather indirectly.

Here agent state & environment state, so we are supposed to build our own agent state

This is called a Partially Observable Markov decision Process

An Agent must construct its own state representation

s_t^a e.g.

complete History : $s_t^a : H_t$

→ Beliefs of the environment state :

$$S^a = \{P[S_t^e = s^1], \dots, P[S_t^e = s^n]\}$$

We build beliefs over where we think we are in the environment

→ Recurrent Neural Network : $S_t^a = \sigma(S_t^a W_s + O_t W_o)$

ii) Components of an RL Agent

Policy :

How an agent goes from a state to picking an action

Value :

How good is each state or action, how much reward do we get if we take a particular action in a particular state

Model :

How agent thinks the environment works, or this is simply the agent's view of the environment

These 3 are not always required, but they may or may not be used

I) Policy is the Agents behaviour

It is a mapping from the agents state to action
The policies could be:

→ Deterministic : $a = \pi_1(s)$

The goal is to make sure the policy picks such an action that maximizes future reward

→ Stochastic Policy : $\pi(a|s) = P[A=a|s=s]$

So this is just the probability of taking a particular action, after being conditioned off a particular state

ii) Value function is a prediction of expected future reward.

It evaluates the goodness/badness of states
(choosing between action/states is done on comparing the total reward one might get between the two.)

$$V_{\pi}(s) = E_{\pi} [R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \dots | s_t = s]$$

Value function depends on which way the agent behaves or rather which action it takes given a state, given a policy

Here γ is the discounting factor based on how valuable the future/distant rewards are.

k) Model Predicts what the environment will do next
 It tries to imagine what the environment does next & then tries to learn the behaviour of the environment, & use that to make a plan & use that to figure out what to do next

We have two parts to our model

→ Transitions : P predict the next state (i.e dynamics)

$$P^a = P[S' = s' | S = s, A = a]$$

Given an action a & current state = s what is the probability we transition to state s'

→ Reward : R predict the next immediate reward

$$R_s^a = E[R | S = s, A = a]$$

The immediate reward we get for taking action a in state 1

This is optional & a lot of agents maybe model free