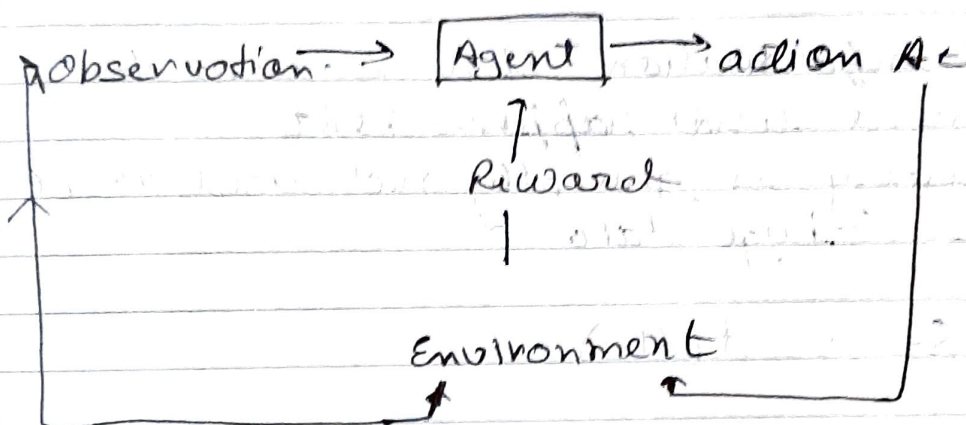


# David Silver

## Lecture 1 →

- Goal of RL → Select actions to maximize total future Rewards  
↳ Series of action to get maximum total result/reward.



Agent Executes Action  $A_t$   
Receives observation  $O_t$   
" scalar reward  $R_t$

Rewards can be scalar but ultimately we want to compare the rewards we can get upon taking different actions so it is preferred to be scalar.

— x — x — x —

→ History → History is the sequence of actions, rewards and observations seen by the agent so far.

$H_t = A_1, O_1, R_1, A_2, \dots, A_t, O_t, R_t$

- ↳ all observable variables upto time  $t$
- ↳ what happens next depends on history

→ State → Information that is used to determine what happens next

- ↳ History is generally not used as it has huge data

$$S_t = f(H_t)$$

This information can be anything → 3 types:-

- ↳ Environment State → The information that is used within the environment. The agent can not see it. It only observes the what happens due to or by this state. For eg. Intermolecular collisions - We don't see them but we know there is water or air or solid on the table.
- ↳ Agent State → The information in the agent's internal representation. The useful observations by the agent which can be used in the algorithm.





~~the~~ ~~re~~

Agent state depends on the function 'f' in

$$S_t^a = f(H_t) \text{ which is decided by us.}$$

↳ Information State  $\Rightarrow$  a.k.a Markov State contains all useful information from the history.

Markov state  $S_t$  is an only if

$$P[S_{t+1} | S_t] = P[S_{t+1} | S_1, S_2, \dots, S_t]$$

the probability

of the next state =

given current

state

The probability of the next

state given all the states

visited so far.

future is independent of the past given the present.

→ × ————— × —————

→ Environment.

↳ Fully observable Environment  $\rightarrow$  Fully able to see what going on in the environment or directly observes the environment.

Agent state = Environment state = Information state

$$O_t = S_t^a = S_t^e$$

Formally this is a MDP (Markov Decision Process)

↳ Partially observable Environment: Agent indirectly observes the environment.

↳ Ex → A blind person trying to know the shape of an ~~obj~~ object. It cannot see it so it indirectly gathers information by touching.

agent state  $\neq$  environment state

→ Partially observable MDP.

- 
- Policy( $\pi$ ) → How the agent picks an action
  - Value function → How good is each state/action
  - Model → Agent's "understanding" of the environment.

→ map from state to action → given a state, it will generate an action.  
can be deterministic or stochastic

↓

$$\pi(a|s) = P(A=a|S=s)$$

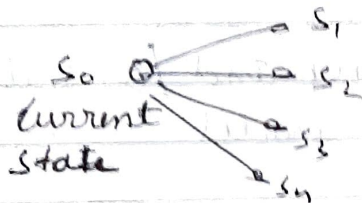
Value Function → It is the prediction of the expected total future reward.



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

↓  
Value function  
for a policy

- model → A model predicts what environment will do next. Its uncertainty.
- transitions → P predicts the next state



each state has some  $p$  of occurring after  $s_0$ . transitions calculate these  $p$ .

- Rewards →  $R$  predicts the immediate reward. It is stochastic in nature.

$$P_{ss'}^a = P(s'=s' \mid s=s, A=a)$$

$$R_s^a = E[R \mid s=s, A=a]$$

Types of RL agents

- Value Based → No Policy  
→ Value Function | observes the values for each action and takes action
- Policy Based → Policy  
→ No Value Function | takes action according to the policy





Date \_\_\_\_\_

Page \_\_\_\_\_

→ Model Free → Policy and/or Value function  
→ No model

We only look at the policy and/or the value function for the particular state. We do not try to understand the working of the environment or make any sort of dynamics of the envt.

→ Model Based → Policy and/or Value Function  
→ No model

We ~~only~~ make an elaborate model of the environment based on the data and try to understand its dynamics

---