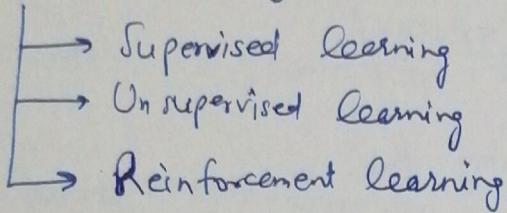


# Introduction to Reinforcement Learning

Lec-1

- Concerned with making a sequence of good / optimal decisions

## Machine Learning



## Characteristics of RL

- There is no supervisor, only a reward signal whether the action was good / bad or +10, -5 etc.
- Feedback / reward is delayed, not instantaneous
- Time really matters (sequential data, not iid)
- In RL, agent ~~takes~~ takes actions that affect the subsequent data it receives

## Examples of RL

- ① Fly stunt manoeuvres in a helicopter
- ② Defeat the world champion at Backgammon
- ③ Manage an investment portfolio
- ④ Make a humanoid robot walk

## Rewards

- \* A reward  $R_t$  is a scalar feedback signal
- \* Indicates how well agent is doing at time  $t$
- \* Job of agent is to sum these  $R_t$  & get as much total reward as possible
- \* RL is based on Reward hypothesis

All goals can be described by the maximization of expected cumulative reward.

## Examples

- ① Fly stunt manoeuvres in a helicopter
  - \* +ve reward for following desired trajectory
  - \* -ve reward for crashing
- ② Defeat the WC at backgammon
  - \* +ve/-ve reward for winning/losing the game
- ③ Manage an investment portfolio
  - \* +ve reward for each \$ in the bank
- ④ Robot walk
  - \* +ve reward for forward motion
  - \* -ve reward for falling over

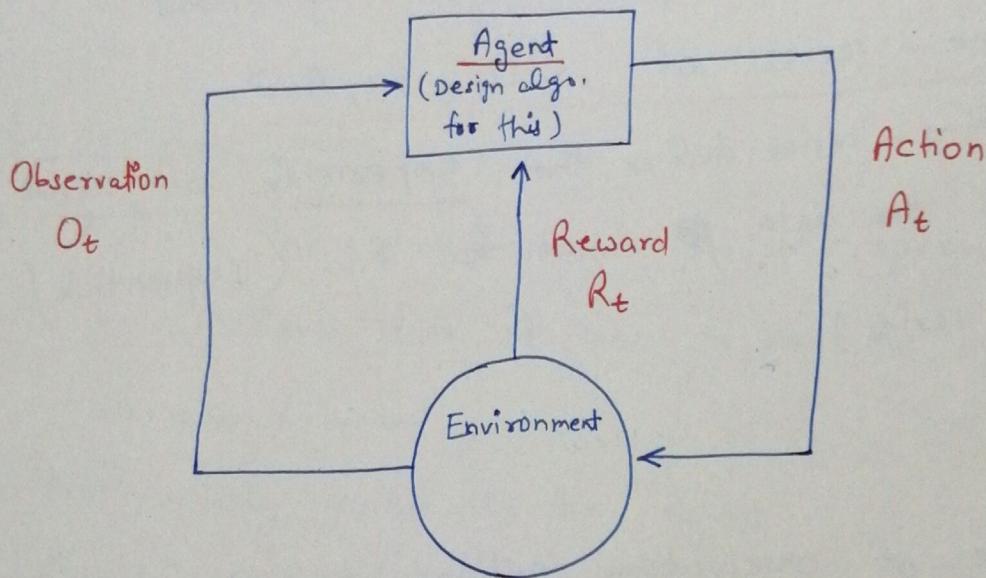
## Sequential decision making

Goal:

Select actions to maximize total future reward

- Actions may have long term consequences
- Reward may be delayed
- It may be better to sacrifice immediate reward to gain more long-term reward

## Agent & Environment



- Environment is the real world out there
- Agent interacts with environment
- We/agent have no control over the environment directly
- Agent influences the environment by actions only

Agent: at each step t

- ① Executes action  $A_t$
- ② Receives observation  $O_t$
- ③ Receives reward (scalar)  $R_t$

Environment: at each time step t

- ① Receives action  $A_t$
- ② Emits observation  $O_t$
- ③ Emits scalar reward  $R_t$

- \* Trial & error loop of agent generates
  - a time-series of  $\langle O_t, R_t, A_t \rangle$
- \* This time-series defines the experience of agent
- \* This is the data used for RL (sequential)  
time-series)

### History

- \* Sequence of observations, actions, rewards

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

i.e. all observable variables up to time t.

- It is all the information that the agent is exposed to
- What happens next depends on the history
  - \* The agent selects actions

\* The environment selects observations / rewards

But, history is often too large / a lot of information to deal with, eg. think of video game.

### State

- Can think it as a summary of history
- State is the information used to determine what happens next
- State is a fun<sup>n</sup> of history :

$$S_t = f(H_t)$$

### Environment state

- $S_t^e$  is the environment's private representation
- i.e. whatever data it uses to pick the next observation / reward
- Not usually visible to the agent
- Even if  $S_t^e$  is visible, it may contain irrelevant info.

### Agent State

- $S_t^a$  is the agent's internal representation
- Info. used by agent to pick next action
- Info used by RL algs
- Can be any fun<sup>n</sup> of history

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

→ Agent can choose fun<sup>n</sup>  $f$  to convert history into state

### Information state (Markov State)

→ Contains all useful info. from the history

→ A state is Markov iff

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

"future is independent of past given the present"

→ If the state representation is Markov, then you can throw away the history

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

→ State is a sufficient statistic of the future

\* The environment state  $S_t^e$  is Markov

\* The history is Markov  $H_t$

→ Agent's state influences the action of agent

→ Different states might / probably will result in different actions

## Fully Observable Environments

- Agent directly observes environment states

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

→ Agent state = Env. state = Info. state

→ formally, this is a Markov Decision Process (MDP)

## Partially Observable Environments

- Agent indirectly observes environment

Eg. Poker playing agent only observes public cards

- Agent state  $\neq$  Env. state

- This is partially observable Markov Decision Process (POMDP)

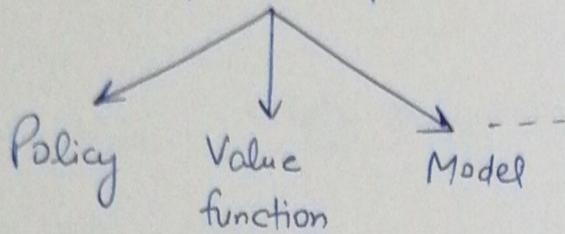
- Agent must construct its own state representation  
 $S_t^a$ . Eg.

① Complete history:  $S_t^a = H_t$

②  $S_t^a = (P(S_t^e=s'), \dots, P(S_t^e=s^n))$

③ RNN:  $S_t^a = \sigma(S_{t-1}^a w_s + O_t w_o)$

## RL Agent components



- Policy is the agent's behaviour
- Map from state to action

Deterministic policy:  $a = \pi(s)$

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

### Value fun<sup>n</sup>

- Prediction of expected future reward
- ~~used~~ used to choose b/w different actions
- Given policy  $\pi$ , its value fun<sup>n</sup>

$$V_\pi(s) = \mathbb{E}_\pi [R_t + r R_{t+1} + r^2 R_{t+2} + \dots | S_t = s]$$

### Model

- A model predicts what the environment will do next

Transition model:  $P$  predicts the next state (i.e., dynamics)

Reward model:  $R$  predicts the next (immediate) reward

$$P_{ss'}^a = P[S'=s' | S=s, A=a]$$

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

In RL agent typically doesn't need all these components to work.

① (a) Value-based agents

- Has value fun<sup>n</sup>
- No policy required
- Take actions greedily based on value fun<sup>n</sup>

(b) Policy-based agents

- Has policy
- No value fun<sup>n</sup>

(c) Actor Critic

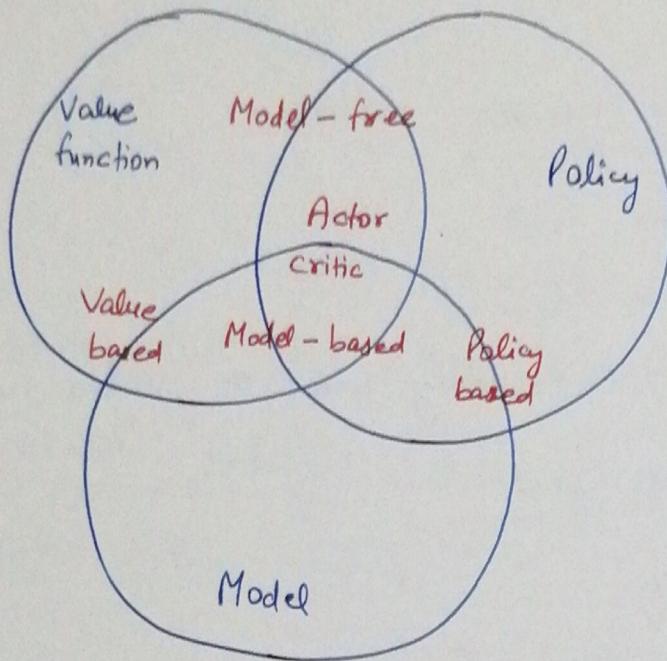
- Policy
- ~~Value~~ Value fun<sup>n</sup>

② (a) Model free agents (\*)

- Don't try to explicitly build model
- Use policy and/or value fun<sup>n</sup>

(b) Model based RL

- Model
- Policy and/or value fun<sup>n</sup>



### Planning

- A model of env. is known
- The agent performs computations with its model (w/o external interactions)
- Agent improves its policy

### Exploration

finds out more info about the environment

### Exploitation

Exploits known information to maximize rewards

- RL is like trial & error learning
- Discover good policy from experiences of the env.  
w/o losing too much reward along the way