

RL (Silver)

Lecture 1.

- RL moves across various fields of science like CS, Neuroscience, Engineering, Psychology, etc.

The reward system in RL algorithms is drawn from the findings in Neuroscience about how the human brain makes decisions, by devoting to the dopamine system.

Machine Learning subdivision

- supervised
- unsupervised
- Reinforcement

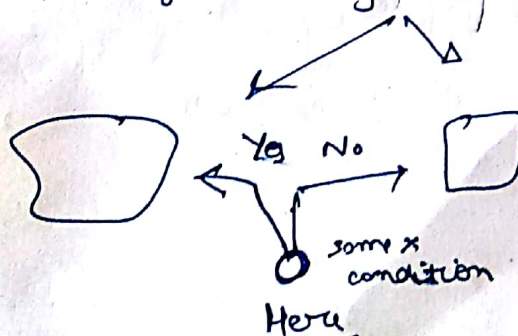
Reinforcement Vs supervised

- No best ~~measured~~ decision mentioned in RL. We just say that for some X action, you get Y points. Not that this is the best action.
- Feedback delayed, not instantaneous \rightarrow In supervised, we directly get the classification result \rightarrow good/bad. But in RL, this decision is delayed, as in, we might get some good points in some of the actions until a catastrophic loss in the end.

(large
-ve
reward)

- non-iid data \rightarrow what we sees in one second is very correlated to what it sees in the next second. So, time really matters.

- Agent gets to make decisions / take actions and based upon its action, it gets very different data.



The RL problem

• Rewards

→ Reward $R \rightarrow$ scalar feedback and there's this hypothesis that any goal can be achieved by maximizing the cumulative rewards at each time steps.

eg. say want to do some task in shortest time then add reward = -1 at each time step, so ^{when} we take the cumulative reward, it learns to take the minimum # steps to maximize cumulative reward.

eg. Defeat world champion at Backgammon.

→ no intermediate reward but final +1/-1 for win/lose.

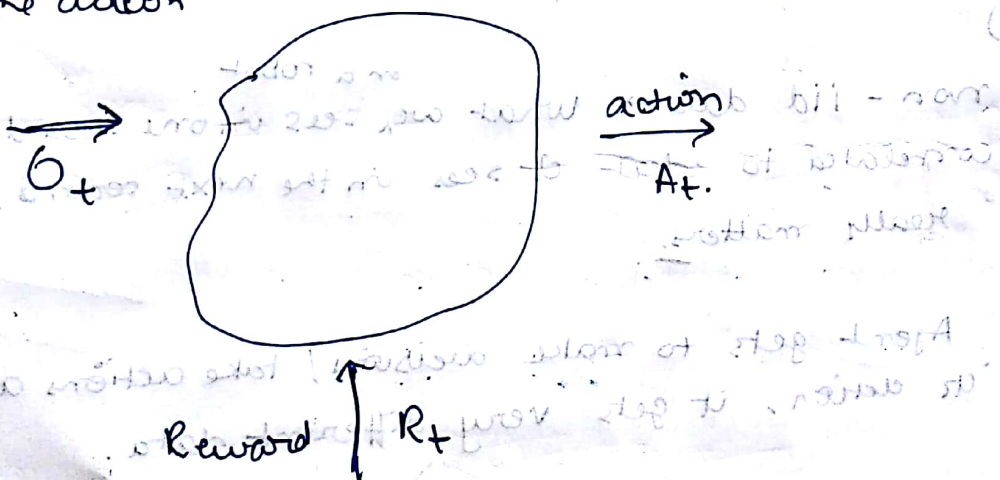
• Sequential decision making

Goal → select actions to max. future total ahead, it may mean to plan ahead and do ~~some~~ short-term losses for long term gains.

• Agent & Environment

↓
responsible for taking the action

Goal → want to build such algo. that sits ~~inside~~ inside the ~~brain~~ agent.



$O_t \rightarrow$ Observations \rightarrow sees something of the real world

$R_t \rightarrow$ Reward \rightarrow how well it's doing at that step.

Environment \rightarrow the world or the thing that the agent interacts with
 \downarrow
responsible for generating O_t & R_t

The only influence the agent can have over the environment is through its action, eg, where it is in the environment by \rightarrow moving to


This goes on in a trial-and-error loop where the agent performs (tries) different actions and based upon its observations & reward, it learns what the most optimal actions are. This time series of O_t , R_t & $A_t \rightarrow$ experience of the agent and forms the data for the RL problem.

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

Goal of our algorithm:

Build a mapping from this history to the next action that the agent should take.

Problem: RL agents may be living for a entire lifetime & hence going back to this history everytime is not feasible.

Soln:  Define a 'state' \rightarrow summary of the history to determine next action.

$S_t = f(H_t)$ / function of History

Environment State : Informatz used within the environment to determine what happens next \rightarrow i.e. some sort of representation (using numbers, etc.) of the environment to generate the next O_t & R_t .

However, this is not visible to the agent \rightarrow i.e. agent/algorithm \rightarrow not ^{be} dependent on the environment state.
 should

Thus, this doesn't influence ~~our~~ building our algorithm.

Agent State \rightarrow the set of numbers used to capture the summary of the informatz seen by the RL agent, which influences what action it should take next. \rightarrow used for building our algorithm.

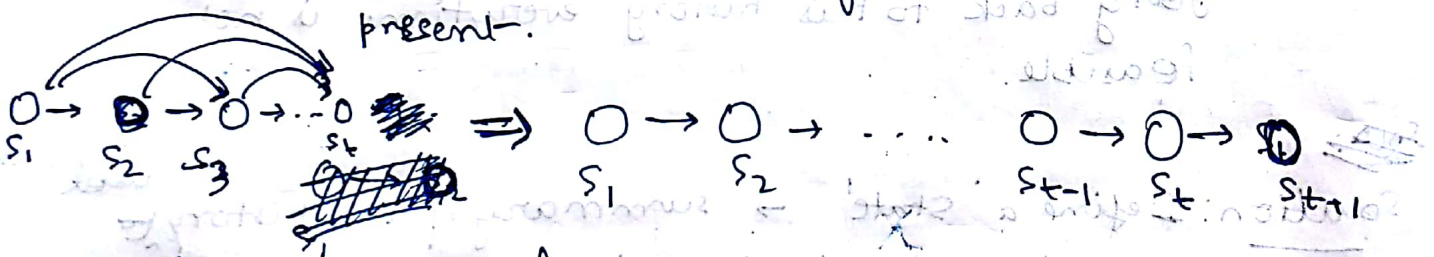
$$S_t^a = f(H_t)$$

Informatz state / markov state \Rightarrow contains all useful informatz from the history.

State S_t is 'markov' if:

$$P(S_{t+1} | S_t) = P(S_{t+1} | S_{t-1}, S_t)$$

i.e. future independent of the past, given the present.



comes from the basic of independency,

than $P(A|B) = P(A)$ if A & B are independent

it is a simplifying assumptz / property.

$$\rightarrow P(A|B, C) = P(A|B)$$

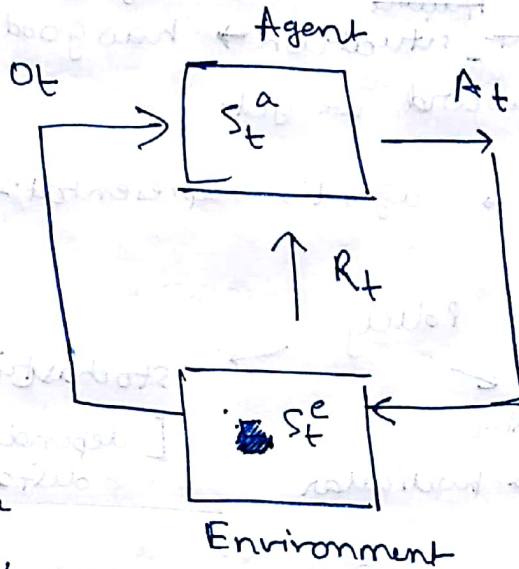
$\rightarrow A$ & C are conditionally independent - given

B.

this means the state $\rightarrow S_t \leftarrow H_{1:t}$ $[H_{1:t} = H_1 H_2 \dots H_t]$
 is sufficient to determine the future $\rightarrow H_{t+1:t}$ \rightarrow determines
 environment state $\rightarrow S_t^e \rightarrow$ markov (by definition)
 we can always find markov states, need to see if we can
 find some that is practical.

Full observability \rightarrow

agent directly
 observes environment
 state $\rightarrow O_t = S_t^e = S_t^a$
 = informativ
 state.



This is the main formula
 behind RL \rightarrow markov
 decision process (MDP).

Partial observability \rightarrow agent indirectly observes the
 environment state.

agent state
 $S_t^a \neq S_t^e$

eg. Robot with camera \rightarrow doesn't know precisely where it
 is \rightarrow needs to figure out by moving & observing

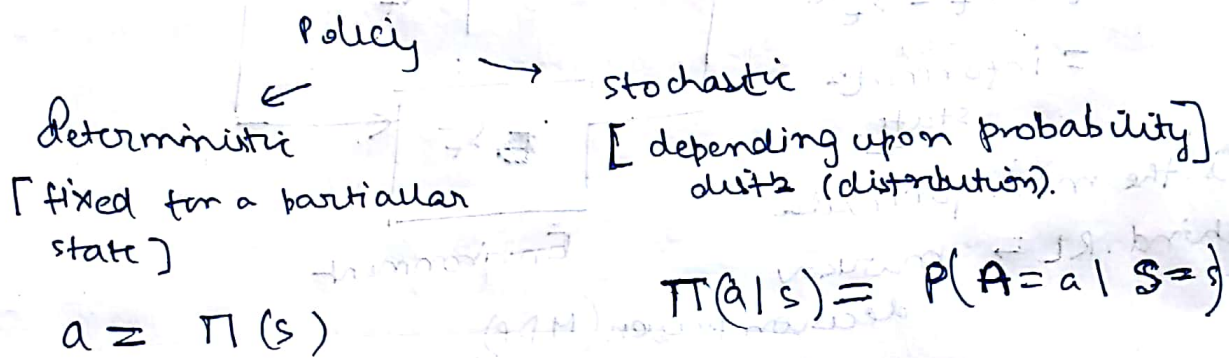
multiple ways to define state:

- ① Naive \rightarrow entire history $\rightarrow S_t^a = H_t$
- ② Probabilistic Bayesian approach \rightarrow probability dists over
 which state it can be in.
- ③ RNNs. (comes in cs224n)

Inside an RL agent

Components (may / may not have)

- Policy → behaviour function → decide from the ^{current} state, which action to take next.
- value function → estimating how well we are doing at the future ~~current~~ ^{future} situation → how good is each state / action, what ~~the~~ reward we get
- model → agent's representation of the environment



Value function → depends on what policy you choose
eg. robot falling down → less reward than robot which keeps moving

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

$$\left[\begin{array}{l} E_X[X] \Rightarrow \text{expected value of } X \Rightarrow \sum_{i=1}^n p_i x_i \\ X \rightarrow \text{can take values } (x_1, \dots, x_n) \\ \text{and } p_i \rightarrow \text{probability of } X = x_i \\ p_i = P(X = x_i) \end{array} \right]$$

~~γ~~ $\gamma \rightarrow$ discounting factor
 $\gamma < 1 \therefore \gamma, \gamma^2, \gamma^3 \rightarrow$ makes the contribution of future rewards subsequently less important than the current reward.

eg. $\gamma = 0.9$

$$R_t + 0.9 R_{t+1} + 0.81 R_{t+2} \dots$$

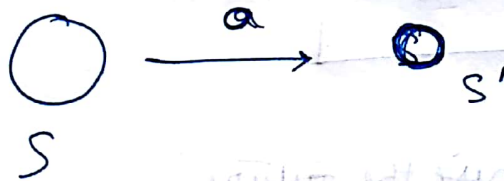
Model \rightarrow since S_{t+1}^a is not visible to the agent, it needs to imagine / predict what environment will do next.

Transitions \rightarrow Rewards \rightarrow predict next immediate reward for action
 $R_{s,a} = E(R | S = s, A = a).$

Given current-state, and action chosen, predict what the next state could be.

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

\uparrow
 State transition probability



RL agent taxonomy:

- \rightarrow value based
 as value function depends on policy, policy is implicit
- \rightarrow Policy Based
 \rightarrow no value function
- \rightarrow Actor Critic
 \rightarrow separate policy
 \rightarrow separate value function.

- \rightarrow model free
- \rightarrow model Based.

2 fundamental problems in RL:

→ General RL

- No environment info known

→ Planning

- model of the environment is given and the agent performs its computations by interacting with this model (not with external environment).

eg. rules of Atari already given and we know what would happen (immediately) if we take a certain step. → Then, we could search over all possible actions and choose the best ones.

[For those who know about AlphaGo, it uses Tree search to find the best possible move]

RL → exploration vs exploitation

Prediction vs Control

- policy given
- predict what will happen if I follow this policy

- optimise the future
- choosing the policy (contingent)
need to solve prediction problem to solve the control problem.