

LECTURE 1:-

* OUTLINE:

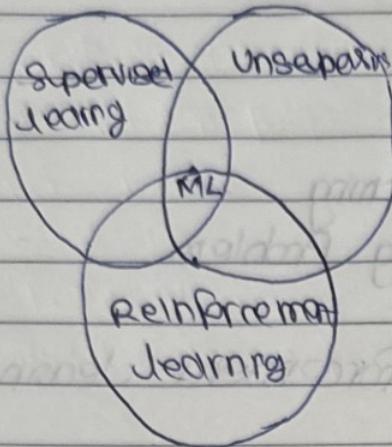
- 1.] Admin
- 2.] About Reinforcement Learning
- 3.] Reinforcement Learning Problem
- 4.] Inside an RL agent
- 5.] Problems within Reinforcement Learning

Book =

- 1.] An introduction to Reinforcement Learning , Sutton and Barto, 1998
- 2.] Algorithms for Reinforcement Learning , Srebro et al

* ABOUT REINFORCEMENT LEARNING:-

RL sits at the intersection of many different fields of science



Supervised learning
Unsupervised learning
Reinforcement learning

Machine Learning

Q.] What makes Reinforcement learning different from other machine learning paradigms?

→ Answer →

- 1.] There is no supervisor, only a reward signal
- 2.] Feedback is delayed, not instantaneous
- 3.] Time really matters (sequential, non i.i.d data)
- 4.] Agent's actions affect the subsequent data it receives

* EXAMPLES OF REINFORCEMENT LEARNING:

- 1] Fly stunt manoeuvres in a helicopter
- 2] Defeat the world champion at Backgammon
- 3] Manage an investment portfolio
- 4] Control a power station
- 5] Make a humanoid robot cook
- 6] Play many different Atari games better than humans

* REINFORCEMENT LEARNING PROBLEM:-

- A reward R_t is a scalar feedback signal
- Indicates how well agent is doing at step t
- The agent's job is to maximise cumulative reward

Reinforcement learning is based on the reward hypothesis

Definition (Reward Hypothesis)

- ⇒ All goals can be described by the maximisation of expected cumulative reward

* Examples of Rewards

- 1] Fly stunt Manoeuvres in a helicopter
the reward for following desired trajectory
the reward for crashing
- 2] Defeat the world champion at Backgammon
the +/- reward for winning/ losing a game
- 3] Manage an investment portfolio
the reward for each \$ in bank
- 4] Control a power station
the reward for producing power
the reward for exceeding safety thresholds
- 5] Make a humanoid robot walk
the reward for forward motion
the reward for falling over
- 6] Play many atari games better than humans
the +/- reward for increasing / decreasing score

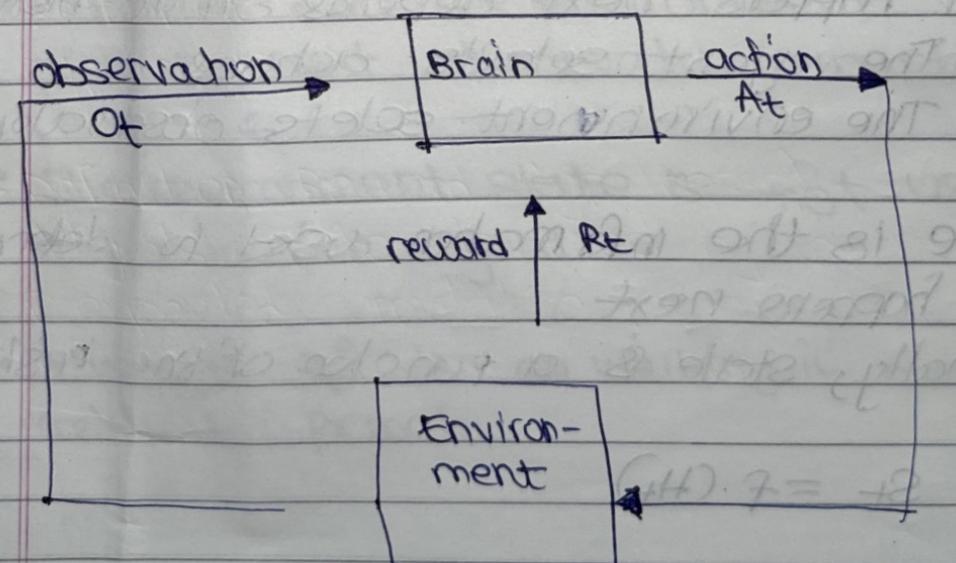
* SEQUENTIAL DECISION MAKING

- Goal:- select actions to maximise total future reward
- Actions may have long term consequences
- Reward maybe delayed
- It may be better to sacrifice immediate reward to gain more long-term reward

Examples:-

- A financial investment (may take months to reach)
- Refuelling a helicopter (might prevent a crash in several hours)
- Blocking opponent moves (might help winning chances many more from now)

* AGENT AND ENVIRONMENT



→ At each step t the agent:

- Executes action a_t
- Receives observation o_t
- Receives scalar reward r_t

→ The environment:

- receives action a_t
- emits observation o_t
- Emits scalar reward r_t

* HISTORY AND STATE:

→ The history is the sequence of observations, actions, rewards

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

→ that is all observable variables up to time t
 → that is the sensor/motor stream of a robot or embodied agent

→ what happens next depends on the history:

→ The agent selects actions

→ The environment selects observations

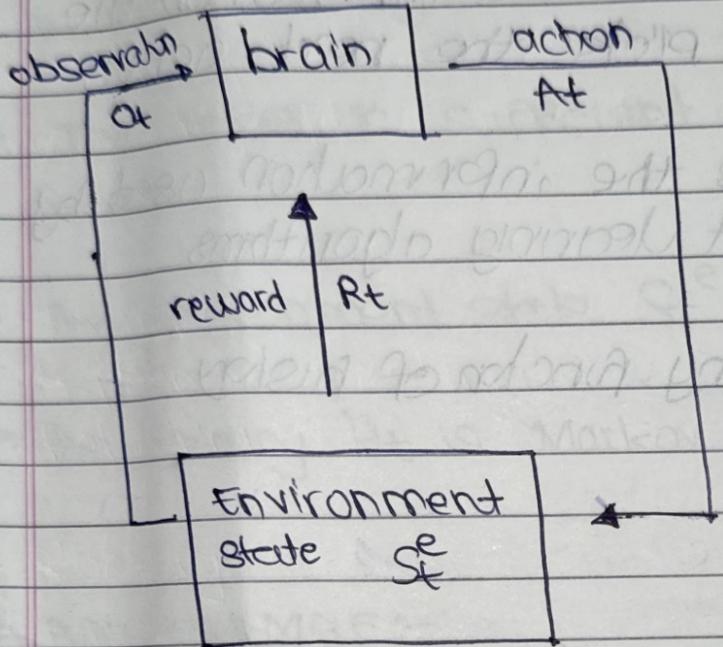
→ state is the information used to determine what happens next

→ Formally, state is a function of the history

$$S_t = f(H_t)$$

BRAIN

* ENVIRONMENT STATE



- The environment state S_t^e is the environment's private representation
- whatever data the environment uses to pick the next observation/reward
- The environment state is not usually visible to the agent
- Even if S_t^e is visible, it may contain irrelevant information

* AGENT STATE:-

- The agent state s_t^a is the agents internal representation
- that is ~~that is~~ whatever information the agent uses to pick the next action
- that is it is the information used by reinforcement learning algorithms
- It can be any function of history

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

* INFORMATION STATE:-

- An information state (a.k.a Markov State) contains all useful information from the history

Definition

A state s_t is Markov if and only if

$$P[s_{t+1} | s_t] = P[s_{t+1} | s_1, s_2, \dots, s_t]$$

↑ probability

- "The future is independent of the past given the present"
- $H_{1:t} \rightarrow S_t \rightarrow H_{t+1:\infty}$
- Once the state is known, the history may be thrown away
- The state is sufficient statistic of the future
- The environment state S_t^e is Markov
- The history H_t is Markov

RAT EXAMPLE:-

light light
 feel feel lever bell electricity
 bell light lever lever cheese
 lever light lever bell

- what if agent state = last 3 items in sequence
- what if agent state = counts for lights, bells, levers
- what if agent state = complete sequence

* FULLY OBSERVABLE ENVIRONMENTS :-

- Full observability :- agent directly observes environment state
 $s_t = s_t^q = s_t^e$
- Agent state = environment state = information
- Formally, this is a Markov decision process (MDP)

* PARTIALLY OBSERVABLE ENVIRONMENTS :-

- Partial observability :- agent indirectly observes environment
 - A robot with camera vision isn't told its absolute location
 - A trading agent only observes current prices
 - A poker playing agent only observes public cards
- Now agent state \neq environment state
- Formally, this is a partially observable Markov decision process (POMDP)
- Agent must construct its own state representation
 s_t^q , e.g.
- Complete history : $s_t^q = h_t$
- Beliefs of environment state $s_t^q = (p_{s_t^q} [s_t^e = s_t^q], p_{s_t^q} [s_t^e = s_t^q], \dots)$
- Recurrent neural network : $s_t^q = \phi(s_{t-1}^q, w_s + \alpha w_o)$

Recurrent neural network - $s_t^a = \sigma(s_{t-1}^a h_t + \alpha)$

* INSIDE AN RL AGENT:-

→ An RL agent may include one or more of these components:

- 1] Policy:- agent's behaviour function and/or
- 2] Value function:- how good is each state action
- 3] Model:- agent's representation of the environment

1] Policy:-

- It is the agents behaviour
- It is a map from state to action, eg
- Deterministic policy : $a = \pi(s)$
- Stochastic policy : $\pi(a|s) = P[A=a | S=s]$

2] Value function:-

- Value function is a prediction of future reward
- Used to evaluate the goodness / badness of states
- And therefore to select between actions, eg

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

3] * model:-

- Model predicts what the environment will do
- Model predicts what the next state will be (ie dynamic)
- Predicts the next immediate reward

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

$$S = S_0, A = a$$

$$R_s^a = E [R | S = S_0, A = a]$$

→ Maze example ✓ ~~policy~~

* CATEGORIZED RL AGENTS

→ Value Based

- No Policy (Implicit)
- Value function

→ Policy Based

- Policy
- No value function

→ Actor Critic

- Policy
- Value function

→ Model free

- Policy and /or Value function
- No model

→ Model Based

- Policy and /or Value function
- Model

⇒ RL agent Taxonomy (diag)

* PROBLEMS WITHIN RL

Two fundamental problems in sequential decision making

→ Reinforcement learning:-

- The environment is initially unknown
- The agent interacts with the environment
- The agent improves its policy

→ Planning:-

- The model of the environment is known
- The agent performs computations with its model (without any external interaction)
- The agent improves its policy
- aka deliberation, reasoning, introspection, pondering, thought search

~~LEARNING~~ → Atari Example

- Rules of game are unknown
- learn directly from interactive game play
- Pick actions on joystick, see pixels and scores

~~PLANNING~~ :-

- Rules of the game are known
- Can query emulator
 - perfect model inside agents brain
- If I take action a from state s
 - what would the next state be?
 - what would the score be?
- Plan ahead to find optimal policy

- Reinforcement learning is like trial and error learning
- The agent should discover a good policy from its experiences of the environment
- without losing too much reward along the way

* Exploration & Exploitation

- Exploration finds more information about the environment
- Exploitation exploits known information to maximise the reward
- It is important to explore & exploit

Examples:-

1] Restraunt Selection:-

Exploitation: go to fav restraunt

Exploration: try new restraunt

2] online banner Advertisements:-

Exploitation Shows the most successful advert

Exploration Shows a different advert

3] Oil drilling

Exploitation: Drill at the best known location

Exploration: Drill at new location

4] Game playing:-

Exploit: play the move you believe is best

Explore: play a new move

* Prediction & Control

- Prediction: evaluate the future
 - Given a policy
- Control: optimise the future
 - Find the best policy

LECTURE-2 MARKOV DECISION PROCESS

Date
Page

- 1] Markov Processes
- 2] Markov Reward Processes
- 3] Markov Decision Processes
- 4] Extensions to MDPs

1] MARKOV PROCESS

→ Introduction to MDPs

- Markov decision processes formally describe an environment for Reinforcement learning
- where the environment is fully observable
- that is the current state completely characterises the process
- Almost all RL problems can be formalised as MDPs eg
 - Optimal control primarily deals with continuous
 - Partially observable problems can be converted into MDPs
 - Bandits are MDPs with one state

2] MARKOV PROPERTY:

"The future is independent of the past given the present"

Definition

A state s_t is Markov if & only if
 $P[S_{t+1} | s_t] = P[S_{t+1} | s_1, s_2, \dots, s_t]$

- The state captures all relevant information from the history
- Once the state is known, the history may be thrown away
- That is State is a sufficient statistic of future

- * Why do we have discount factors?
- To indicate that our model is not perfect

Value Function

eg

If dropped in class '2' from there till pass how many rewards would you get

$$V(S) = E [G_t | S_t = S]$$

Expected return

cg 2

~~$$V(S) = 0.4(0.8) + 0.6(1-0.8) \cdot 2$$~~

~~$$= 0.32 + 1.2 - 2$$~~

$$[e^{-0.05 \cdot 1} / 1.05] = [e^{-0.05 \cdot 2} / 1.05]^2$$

$$V^\pi = R^\pi + \gamma P^\pi V^\pi$$

avg reward
function

avg transition
dynamics

LECTURE 3:- Planning By Dynamic programming

policy evaluation \rightarrow if someone gives you a policy how good is that

policy iteration \rightarrow using policy evaluation

$\pi' \rightarrow$ new policy

→ LECTURE 4 MODEL FREE PREDICTION

Date
Page

- 2] Monte-Carlo → methods that go to the end of trajectory & estimate the value just by looking at sample
- 3] Temporal-Diff
to look 1 step ahead & estimate return

MC / TD better?

TD is more efficient?
because markov property

→ LECTURE 5:- Model free control

when you throw your robot in an unknown environment

TD learning $\xrightarrow{\text{for}}$ off policy