

# RL Lecture 2

## Markov Decision Process

### → Markov Process:

Intro to MDPs:

- MDPs formally describe the environment for RL.
- Environment is fully observable.  $O_t = S_t^e = S_t^a$ .
- Almost all RL problems → can be formalized as MDPs.

Owing to the Markov property,

the state captures all relevant information from the history.

State transition matrix:

$$P_{ss'} \equiv P[S_{t+1} = s' \mid S_t = s]$$

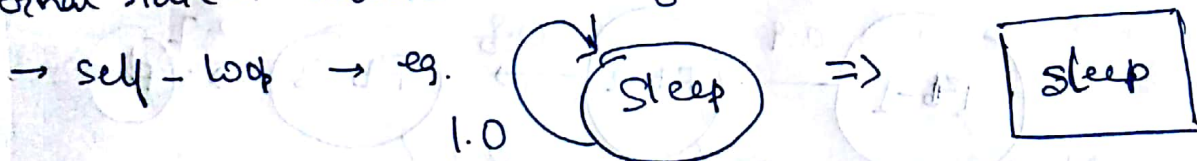
$P \rightarrow$  transition probability  $\rightarrow$

$$\begin{bmatrix} P_{11} & \dots & P_{1n} \\ \vdots & & \vdots \\ P_{n1} & \dots & P_{nn} \end{bmatrix}$$

Markov process:

Random process (sequences of states  $s_1, s_2, \dots$ ) following the Markov property.

Internal state  $\rightarrow$  doesn't need any modelling



~~Modifiers:~~  
Modifiers:

$\rightarrow$  probability is changing over time.

2 ways modelling

① model it as a non-stationary process.

stationary process:  
For an  $\infty$  sequence of states, it is called a stationary process when:

$$P(X_1, X_2, \dots, X_t) = P(X_{1+k}, X_{2+k}, \dots, X_{t+k})$$

basically

$$\dots (X_1, X_2, \dots, X_t)$$

$$\dots (X_{1+k}, \dots, X_{t+k}) \dots$$

$t$  consecutive at any point of time.

where  $P(X_1, X_2, \dots, X_t) \rightarrow$  joint distribution]

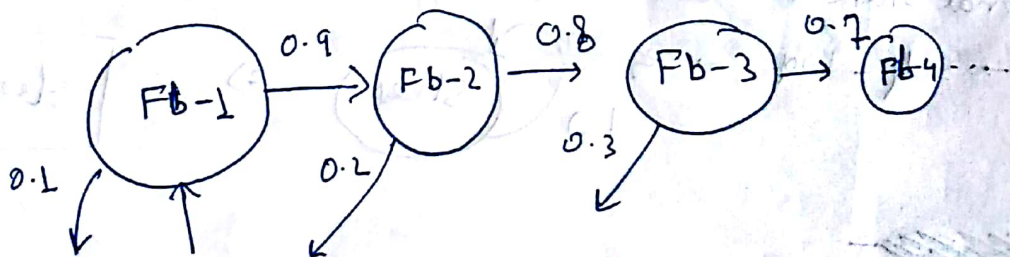
So, the 1st proposed sol<sup>n</sup> is to use non-stationary/  
MDP and incrementally adjust your sol<sup>n</sup>.  
(later)

② non-stationary dynamics  $\Rightarrow$  just a more complicated Markov process

eg.

Facebook

in the simple case  
Becomes



Fb- $i$  represents that you spent  $i$  seconds in that job (here,  $i$  seconds on facebook).



## Markov Reward Process

→ represented by tuple  $(S, P, R, \gamma)$

→ Markov chain with values

We care about the max of this reward i.e. maximizing this reward.

$S$  → set of states

$P$  → transition probability

$R$  → Reward function

(immediate)

$= E(R_{t+1} | S_t = s)$

$\gamma$  → discount factor

$$G_T = R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} \dots$$

(return)

to make sure that it is finite.

$$= \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

Note: No expectation here, as

$G_T$  → is just a sample, not the overall value. at a time instant.

$\gamma = 1$  → maximally short / far sighted  
Why?

→ mathematically convenient

→ model is not perfect → uncertainty of the future

$G_T$  → means from time-step onwards.

Value function → long-term value of being in a state.

$$v(s) = E[G_T | S_T = s]$$

$R_{t+1}$  → indexing doubt → just a convention

→ stick to the concept that:

- move to state
- get a reward.

Bellman equation for MRP:

Value function → two parts

$R_{t+1}$

discounted value of successor state →  $\gamma v(t+1)$

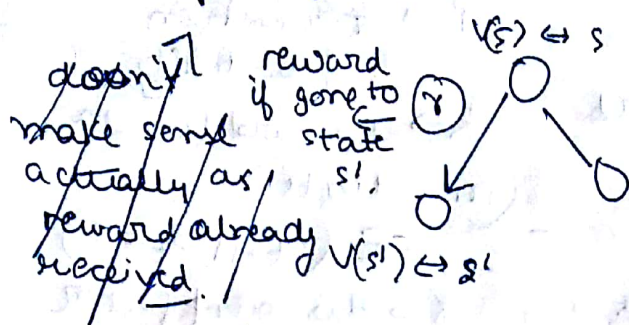
$$v(s) = E[G_T | S_T = s]$$

$$= E[R_{t+1} + \gamma G_{t+1} | S_t = s] = E[R_{t+1} + \gamma v(t+1) | S_t = s]$$



Value function needs to obey Bellman equations

Backup diagram  $\leftrightarrow$  one step look ahead search



$$V(s) = R_s + \gamma \sum_{s' \in S} P_{ss'} V(s')$$

matrix formulation:

$$V = R + \gamma P V \rightarrow O(n^3) \rightarrow \text{not feasible to compute directly}$$

$$V = (I - \gamma P)^{-1} R$$

$\rightarrow$  use DP / Monte-Carlo / Temporal ~~Difference~~ ~~Process~~ Difference Learning

Then now, no actions

Enter:  $\downarrow$  MDP  $\rightarrow$  MRP with decisions (or actions).

Environment in which all states are Markov.

add  $A$  to tuple  $\rightarrow$  set of actions.

$$P_{ss'}^a = P(S_{t+1}=s' | S_t=s, A_t=a)$$

$\downarrow$  dependent on 'a' now

$$R_s^a = E[R_{t+1} | S_t=s, A_t=a]$$

We need to formalize the meaning of making decisions and taking decisions.

$\downarrow$  Policy  $\rightarrow$  distr over ~~set~~ actions given states.

$$\pi(a|s) = P(A_t=a | S_t=s)$$

$\rightarrow$  markov policy  $\rightarrow$  only dependent on current state & not the time step you are in.



policy  $\rightarrow$  fully characterizes the behaviour of an agent.  
 $\hookrightarrow$  same at each time instant  $\rightarrow$  just depends on the state.

[future rewards  $\rightarrow$  fully characterized by state].

We can always recover a MRP from a given MDP.

No matter what policy chosen, it defines some chain of ~~sequences~~ <sup>states</sup> which ultimately form a Markov chain

& sequences of states & rewards  $\Rightarrow$  MRP  $\langle S, P^\pi, R^\pi, \gamma \rangle$

$\rightarrow$  Just average over the policy.

$$P_{ss'}^\pi = \sum_{a \in A} \pi(a|s) P_{ss'}^a$$

$$R_s^\pi = \sum_{a \in A} \pi(a|s) R_s^a$$

Earlier, value functions didn't have any action associated with it. Now:

① state value function:  $V_\pi(s) \rightarrow$  expected return starting from state  $s$  and following policy  $\pi \rightarrow$  sampling actions according to  $\pi$ .

$$V_\pi(s) = E_\pi(G_T | S_T = s)$$

② action-value function  $\rightarrow$  expected return starting from  $[q_\pi(s, a)]$  state  $s$ , taking action  $a$  & following policy  $\pi$ .

$$q_\pi(s, a) = E_\pi[G_T | S_t = s, A_t = a]$$

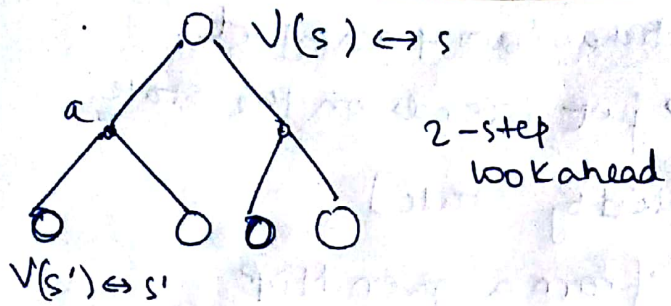
[Difference: here action at this state is chosen as  $a$  and then find the expectation].

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

$$q_\pi(s, a) = E_\pi[R_{t+1} + \gamma q_\pi(S_{t+1}, A_{t+1}) | S_t = s, A_t = a]$$

$$V_\pi(s) = \sum_{a \in A} \pi(a|s) q_\pi(s, a)$$

$$q_\pi(s, a) = R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V_\pi(s')$$



$$V_{\pi}(s) = \sum_a \pi(a/s) \left[ R_s^a + \gamma \sum_{s'} P_{ss'}^a V_{\pi}(s') \right]$$

$$q_{\pi}(s, a) = R_s^a + \gamma \sum_{s'} P_{ss'}^a \sum_{a'} \pi(a'/s') q_{\pi}(s', a')$$

↳ Recursive relation.

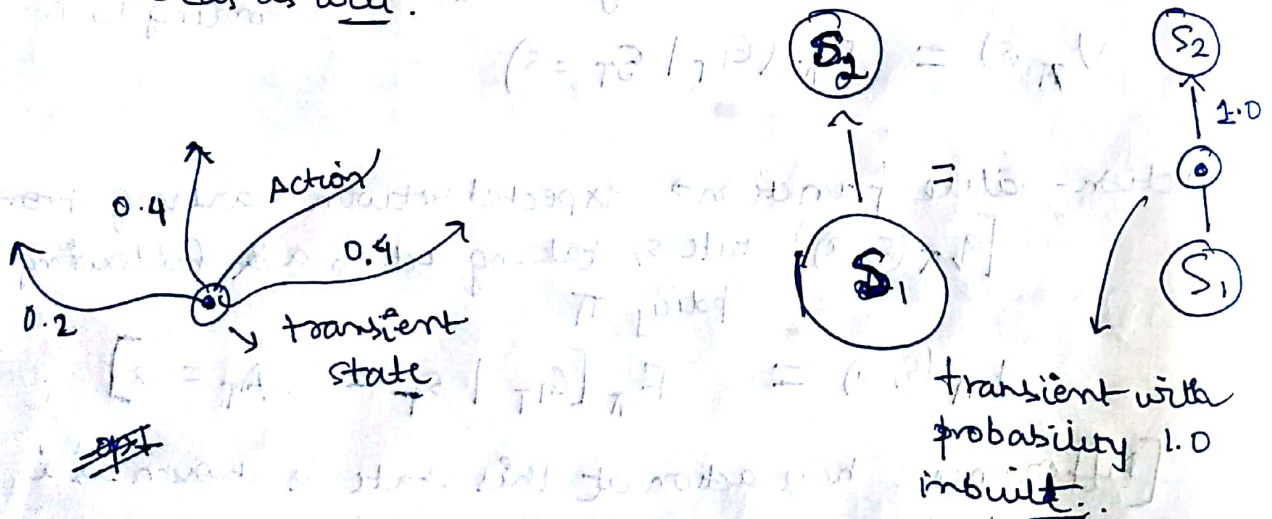
1 way to solve → use the matrix formulation:

$$V_{\pi} = R^{\pi} + \gamma P^{\pi} V_{\pi} \quad (\text{framing MDP as a MRP})$$

$$V_{\pi} = (I - \gamma P^{\pi})^{-1} R^{\pi} \rightarrow \text{more efficient methods later.}$$

But how to get the optimal?

Reward → dependent on action →  $R_s^a$   
 i.e. actions can decide what happens with the immediate rewards as well.





Optimal value function :  $V_*(s) = \max_{\pi} V_{\pi}(s)$

$q_*(s, a) = \max_{\pi} q_{\pi}(s, a) \rightarrow$  max, following ~~the~~ all the different kinds of policy possible

Knowing the optimal value function

$\Rightarrow$  MDP solved!!

Optimal policy

$\leftarrow$  best possible way to behave

mapping from state to action

which is best?

partial ordering over policies:

what is optimal?

$$\pi \geq \pi' \text{ if } V_{\pi}(s) \geq V_{\pi'}(s) \forall s$$

Imp theorem:

$\rightarrow \exists$  atleast 1 optimal policy  $\pi_* \geq \pi \forall \pi$

$\rightarrow$  can be more than 1 optimal policies  $\rightarrow$  all with same value function:  $V_{\pi_*}(s) = V_*(s)$

$\rightarrow$  they also get the optimal action-value function:

$$q_{\pi_*}(s) = q_*(s)$$

Finding  $\rightarrow$

$$\pi_*(a|s) = \begin{cases} 1 & \text{if } a = \max_{a \in A} q_*(s, a) \\ 0 & \text{else.} \end{cases}$$

max over

~~$q_{\pi}(s, a)$~~   $q_*(s, a) \rightarrow$  solve for  $q_*$   $\rightarrow$  done  $\checkmark$

$\rightarrow$  ~~also~~ always a deterministic optimal policy

How to get  ~~$q_{\pi}(s, a)$~~   $q_*(s, a)$ ?

Before, we looked at Bellman expectation equation.

Bellman optimal equation

$\rightarrow$  1 step look ahead

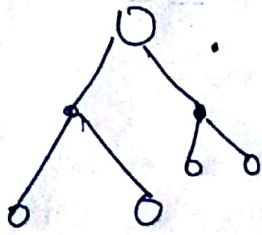
optimal value functions  $\rightarrow$  recursively related.

$$V_*(s) = \max_{a \in A} q_*(s, a)$$

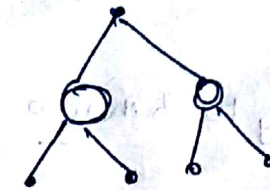
$$q_*(s, a) = R_s^a + \gamma \sum_{s'} P_{ss'}^a V_*(s')$$

$$V_{\pi}(s) = \max_a (R_s^a + \gamma \sum_{s'} p_{ss'}^a V_{\pi}(s'))$$

$$Q_{\pi}(s, a) = R_s^a + \gamma \sum_{s'} p_{ss'}^a \max_{a'} Q_{\pi}(s', a')$$



(Risk sensitive MDPs?)



makes eq 2  
non-linear.  
Hence, we  
iterative methods

value iteration,  
policy iteration,  
Q-learning etc...