

Lecture 2

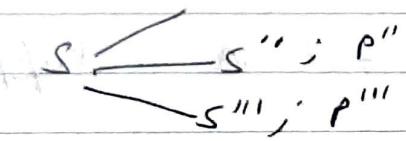
1) A state s_t is markov if and only if

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

we dont need to know s_1 to s_{t-1} to understand $P[s_{t+1}]$. s_t is sufficient if it is an MDP

→ State Transition matrix

the transition probability from state s to s' in an episode is $P_{ss'} = P[s_{t+1} = s' | s_t = s]$
 if we wait, then $P_{ss''} = P[s_{t+1} = s'' | s_t = s]$



State transition matrix is just a matrix representation of these probabilities.
 i.e. $P_{ss'}$ or $P_{ss''}$...

We are in state s , what is the probability of us going to state s' from there $\rightarrow P_{ss'}$

For n states

$$P =$$

$$\begin{array}{c} \xleftarrow{\quad} \xrightarrow{\quad} \\ \downarrow \\ \int \begin{bmatrix} P_{11} & \cdots & P_{1n} \\ \vdots & & \vdots \\ P_{n1} & \cdots & P_{nn} \end{bmatrix} \end{array} \quad \begin{array}{l} \text{sum of each} \\ \text{row is 1} \end{array}$$

each ~~row~~ element represent transition Prob.

P from this state to this state

State transition matrix gives us the complete structure of the markov problem. It tells us the complete picture, from one state (the row number), how likely I am going to the next state (the column number).

→ A markov process is therefore a random process which satisfies markov property

It can be represented as a tuple

$\langle S, P \rangle$

- S is a finite set of State (State Space)
- P is a state transition probability matrix.

It is sort of a representation of the system in which our agent is operating

→ Markov Reward Process \Rightarrow Markov process with values \rightarrow how good is it to be in a particular state

→ It is represented as a tuple: the odd reward and discount factor γ to the previous tuple $\langle s, P, R, \gamma \rangle$

R is a reward that I get from a state at that moment

We want to maximize the cumulated sum of rewards in an episode

$$R_s = E[R_{t+1} | s_t = s]$$

it is stochastic

$\gamma \in [0, 1]$: it is used to regulate the total reward. Whether we want to give more weight to future rewards or the immediate reward.

↳ Returns \rightarrow this is what we care about
 It is the total discounted reward from time step t

$$G_{t+1} = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum \gamma^k R_{t+k+1}$$

Note: $t+1$ because we were in a state at everything after that happens at $t+1$

G_t is the return generated from a random episode observed by the agent out of all the episodes possible

$$\rightarrow \text{Value Function} \rightarrow V(s) = E[R_t | s_t = s]$$

Value function is basically, if we are at some state s , what is the expected value of R_t that we can get given all the possible episodes that can happen after this state and R_t for those episodes. Episodes that have higher probabilities of happening will "contribute" more to the expected value. We have expectation because the rewards are stochastic.

So we have two ~~one~~ stochastic variables here which will determine the value function. 1) Reward. 2) Transition probability.

$$\rightarrow \text{Bellman Equation for MRB}$$

\hookrightarrow Opt

$$\hookrightarrow V(s) = E[R_t | s_t = s]$$

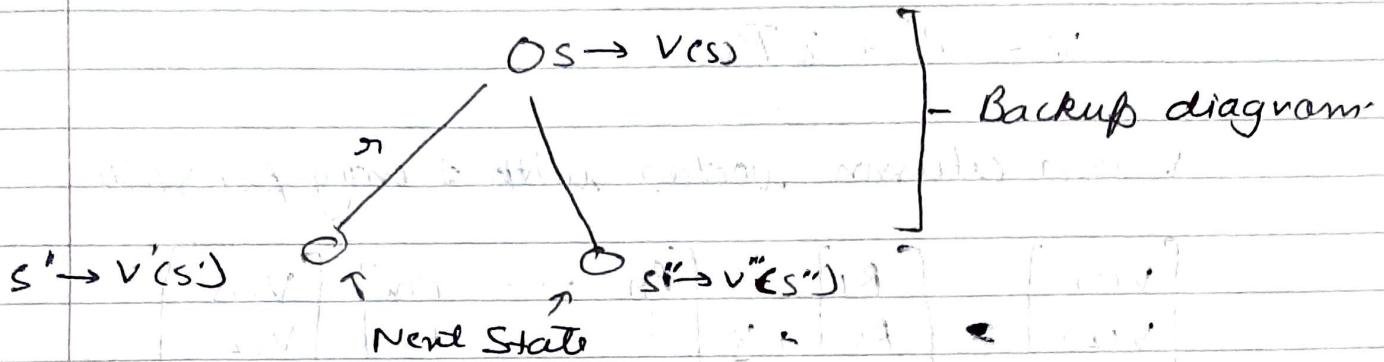
$$V(s) = E[R_{t+1} + \gamma R_{t+2} + \dots | s_t = s]$$

$$V(s) = E[R_{t+1} + \gamma [R_{t+2} + \gamma R_{t+3} + \dots] | s_t = s]$$

$$V(s) = E[R_{t+1} + \gamma V(s_{t+1}) | s_t = s]$$

$$V(s) = E[R_{t+1} + \gamma V(s_{t+1}) | s_t = s]$$

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



We can calculate the value function of state s given we know the value function of all the next states, their transition probability and the reward for state s .

Consider we are at state s and the next possible states are s' and s'' . Calculate the value function of state s

↳ through bellman equation →
Since we have multiple possible states for the next step. To know the expected value of $t+1$ for state s , we will sort of average the $V(s_{t+1})$ using the transition probabilities of all the possible states for the next step.

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

↓ reward for state s ↓ state space.

→ Bellman Equation in Matrix form

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

V is a column vector with 1 entry per state

$$\begin{bmatrix} V(1) \\ V(2) \\ \vdots \\ V(n) \end{bmatrix} = \begin{bmatrix} R_1 \\ R_2 \\ \vdots \\ R_n \end{bmatrix} + \gamma \begin{bmatrix} P_{11} & \cdots & P_{1n} \\ \vdots & \ddots & \vdots \\ P_{n1} & \cdots & P_{nn} \end{bmatrix} \begin{bmatrix} V(1) \\ V(2) \\ \vdots \\ V(n) \end{bmatrix}$$

$$\Rightarrow (I - \gamma P) V = R$$

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

$O(n^3)$ for n states



→ Markov Decision Process \rightarrow Markov reward process with decisions given all the states are markovian

It is a tuple $\langle S, A, P, R, \gamma \rangle$

action set

A is a finite set of action the agent can take

$$P \rightarrow P_{ss'}^a = P[S_{t+1} = s' / S_t = s, A_t = a]$$

state transition probability from s to s' given action a

R is a reward function for state s given action a $R_s^a = E[R_{t+1} / S_t = s, A_t = a]$

P is now dependent on the action. So we can have separate transition probability matrix for each action in the action set $[A, n, n]$

$\times \quad - \quad \times \quad - \quad \times$

→ Policy \rightarrow Stochastic in nature \rightarrow Distribution of action given state \rightarrow Each action has some probability of being taken given the state.

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



as it is an MDP, the policy only depends on the current states

⇒ we have ~~stational~~ stationary policy i.e. the policy for a state is independent of the time step. It remains the same.

In practice, this may not be true. Our environment may not support MDP. For eg. Partially Observed Environment:

Now we have another stochastic variable, π which affects the R for a state and the transition probability.

⇒ the average reward for a state given a policy would be an avg over all actions

$$R_{S,\pi} = \sum_{a \in A} \pi(a|s) R_{S,a}$$

⇒ the average transition probability for all actions given π from s to s' would be an average of $P_{S,S'}$'s over all actions

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

→ Value function is therefore dependant on policy as some actions have more chances of being taken than other and therefore the rewards and the transition probabilities beyond that state

$$\Rightarrow V_{\pi}(s) = E_{\pi} [G_t | S_t = s]$$

How good is it to be in a state s if we follow the policy π .

→ Action Value function \Rightarrow Unlike the state value function which tells us how good is it to be in a particular state, action value function tells us how good is it to take a particular action. It's a criterion to evaluate and compare all the possible action from the given state.

It is the expected return starting from state s , taking action a , following and then following policy π .

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

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

~~$$q_{\pi}(s) = E_{\pi} [R_{t+1} + \gamma q_{\pi}(s_{t+1}) | S_t = s]$$~~

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

$$V_n(s)$$

a

$$q_n(s,a)$$

$$V_n(s) = \sum_{a \in A} \pi(a|s) q_n(s|a)$$

probability of taking that action.

$$q_n(s,a)$$

a

$$q_n(s,a)$$

$$V_n(s')$$

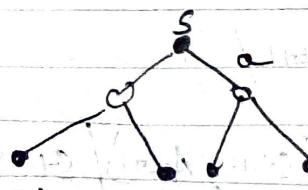
s'

~~Maximize~~

$$q_n(s,a) = R_s^a + \sum_{s' \in S} P_{ss'}^a V_n(s')$$

combining ① and ②

$$V(s) = \sum_{a \in A} \pi(a|s) \left[R_s^a + \sum_{s' \in S} P_{ss'}^a V_n(s') \right]$$



$$\rightarrow \text{Matrix form} \rightarrow V_n = R^\pi + \gamma P^\pi V_n$$

$$V_n = (1 - \gamma P^\pi)^{-1} R^\pi$$

Optimal Value Function -

$V^*(s) \rightarrow$ It is the max value function over all policies \Rightarrow The max possible value of the value function at a state
 \rightarrow The value function that we get for the best policy i.e. the maximum value possible reward

$$V^*(s) = \max_{\pi} V_{\pi}(s)$$

Optimal action value function is the maximum action value function over all policies.

$$q^*(s,a) = \max_{\pi} q_{\pi}(s,a)$$

To behave optimally at state s we pick the action which has the highest $q^*(s,a)$ and therefore for this policy the $V^*(s) = q^*(s,a) + R$



Date _____
Page _____

→ Optimal policy \rightarrow Best possible way to behave in an MDP.

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

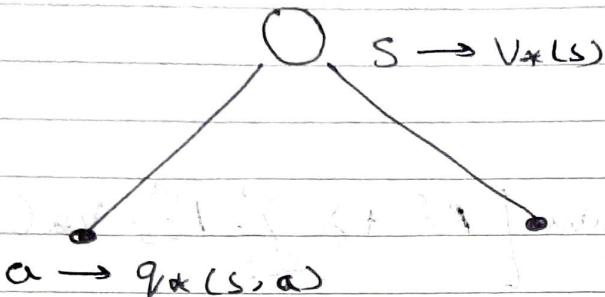
If the value function of π > value function of π' for all states then
 $\pi > \pi'$

For any MDP, there is always an MDP which is better \Rightarrow than any other policy (or equal to)

It is possible to have more than one optimal policy. If they are 2 optimal, they will have the same value function i.e. the optimal value function and also the optimal state-value function.

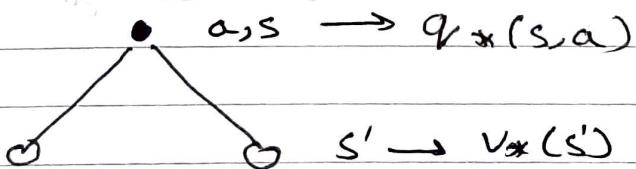
Optimal policy can be found by maximizing over $q_\pi(s, a)$

→ Bellman optimality equations



For optimality, we want to pick the action & having maximum q^* . So instead of taking the average of the q values for all the actions, the v^* would be the max q^* over all the actions

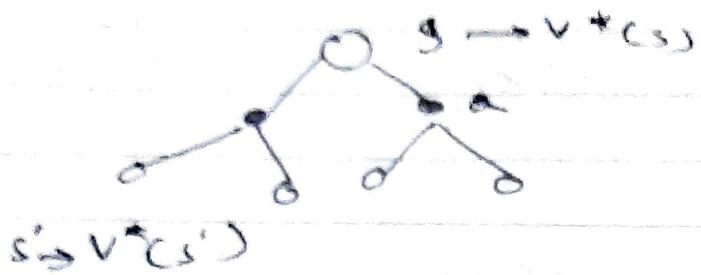
$$v^*(s) = \max_a q^*(s, a)$$



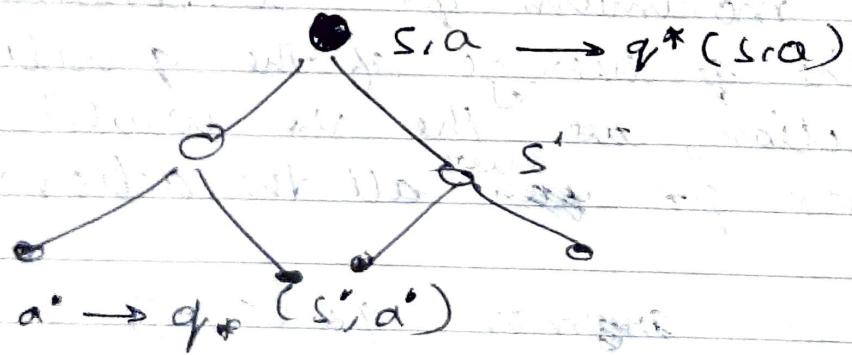
$$q^*(s, a) = r_s^a + \gamma \sum_{s' \in S} p_{ss'}^a v^*(s')$$

We don't get to pick the next state. It happens according to the transition probability so we average the v^* values for the next states according to their transition probabilities.

W



$$V^*(s) = \max_a \left[R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V^*(s') \right]$$



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