

## Markov Decision Process

a) Markov Decision Process formally describes an environment for reinforcement learning, where the environment is fully observable to the agent.

This simply means that the current state fully characterizes the decision process.

### b) Markov Property

A State is Markov if & only if

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

The State  $S_t$  captures all relevant information for the State  $S_{t+1}$  to depend on it that it can discard all the previous states & still have no effect.

Once the state  $S_t$  is known, we can throw away the history.

c) For a Markov State  $s$  & its successor State  $s'$ , we can define a probability that we transfer from State  $s$  to  $s'$  by

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

Since  $S_t = s$  is a markov state it means it

characterizes all information about the history & hence given such a state we have all the information to figure out what will happen next.

So if we imagine  $n$  possible states, for any given state assuming it is markov we can define a matrix with the probability that, given a particular state what is the probability we move into the next state.

$$P = \begin{matrix} & \left[ \begin{matrix} P_{11} & \dots & P_{1n} \\ \vdots & \ddots & \vdots \\ P_{n1} & \dots & P_{nn} \end{matrix} \right] \\ \text{from} & \end{matrix}$$

here each row of the matrix sums upto 1

We can follow this through to multiple processes & continue sampling through next possible states.

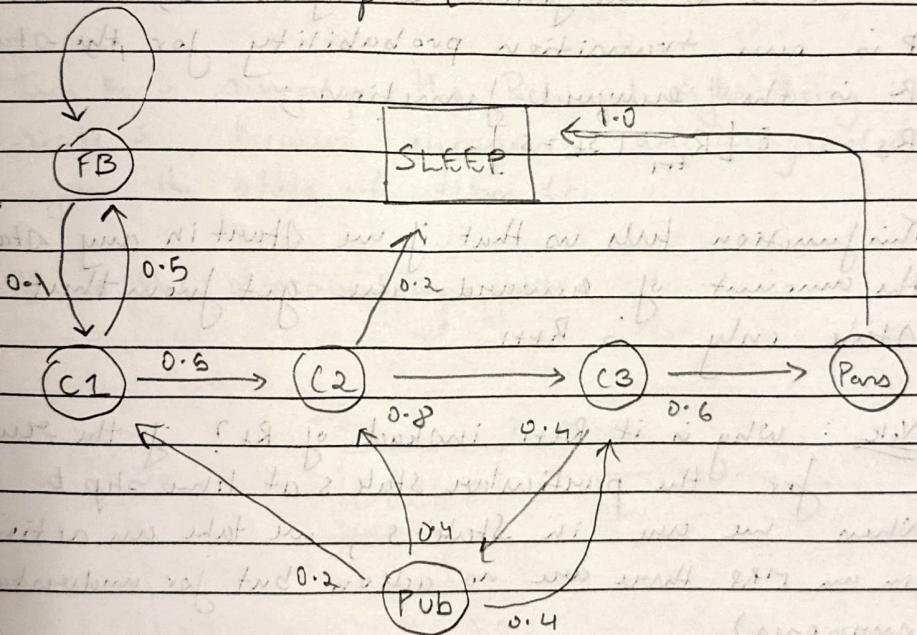
D) A markov Process is hence a memory less random process i.e a sequence of Random States  $S, S_1, S_2, \dots$  with the markov property

So a Markov Process (or a markov chain) is a tuple  $(S, P)$  where  $S$  is the finite set of States (the possibilities of where we can be in) &  $P$  is the state transition probability matrix,

$P$  is a state transition probability matrix, which characterizes how a transition from one state to the next happens. The entries in the matrix are all the probabilities of moving from one state to another.

Example: (i) Find probability of getting bus from school

	Class 1	Class 2	Class 3	Pass	FB	Pub	Sleep
P =	Class 1	0.5	0.5	0.5	0.5	0.5	0.5
	Class 2	—	0.8	0.8	0.8	0.8	0.2
	Class 3	—	0.6	0.6	0.6	0.4	—
	Pass	—	—	—	—	—	1
	Facebook	0.1	—	—	—	0.9	—
	Pub	0.2	0.4	0.4	0.4	0.4	—
0.9	Sleep	—	—	—	—	—	1



(E) Markov Reward Process is a Markov chain with values which assign a reward / value judgement which says how good it is to be in or how much reward is accumulated across a particular sequence sampled from this state transition matrix.

A markov reward process is a tuple  $\langle S, P, R, \gamma \rangle$   
→ where  $S$  is the finite set of states,  
→  $P$  is the transition probability for the states  
→  $R$  is the reward function;  
$$R_S = E [R_{t+1} | S_t = s]$$

This function tells us that if we start in any state  $s$ , the amount of reward we get from that state only is  $R_{t+1}$ .

Note : Why is it  $R_{t+1}$  instead of  $R_t$ ? if the rewards for the particular state  $s$  at time step  $t$  when we are in state  $s$ , we take an action (in an MRP there are no action but for understanding purposes)

The environment moves to state  $S_{t+1}$  & gives reward  $R_{t+1}$ .

So reward belongs to a transition  
Reward for being in a state is same as  
Reward emitted when we transition out of state  $s$

eg In a grid world

we are at state  $s = (3, n)$

this is a five tile (bad location)

so each time we are in this state if we

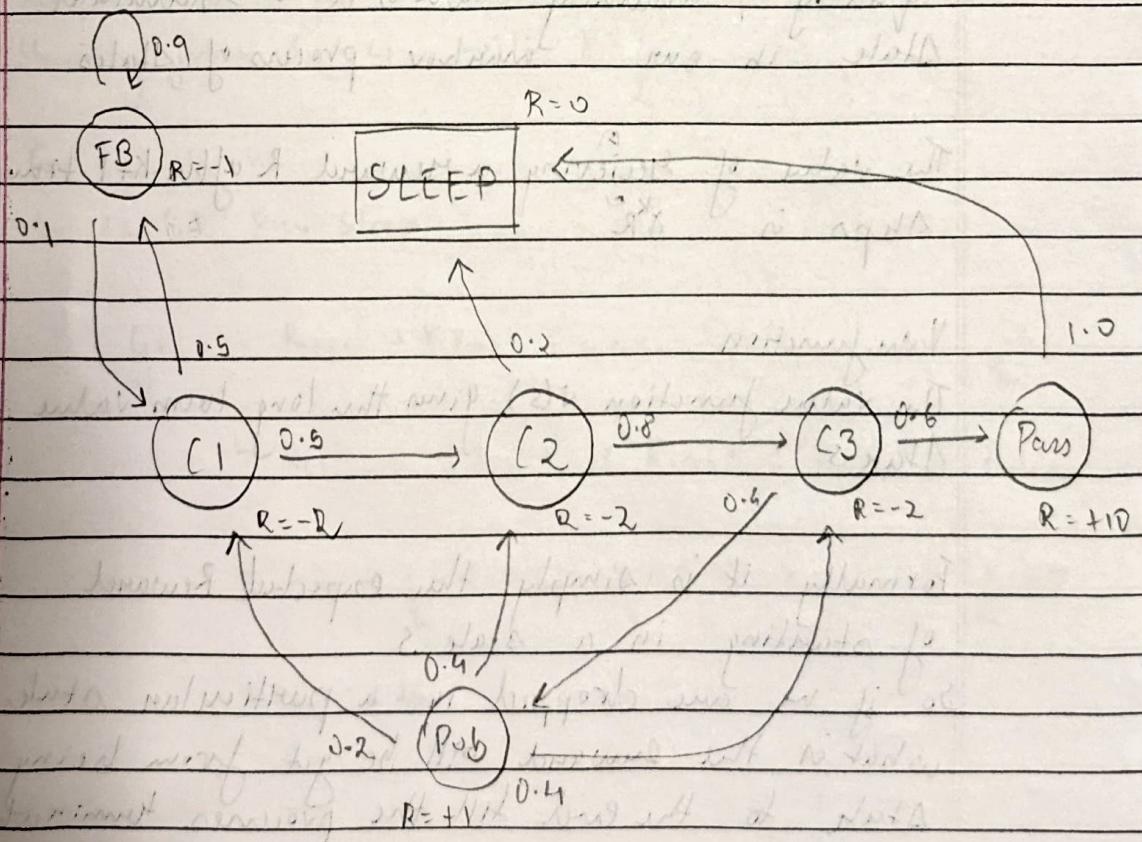
enter of continue being in this state

we receive a -5 reward at the next  
time step.

$$\therefore R_s = E[R_{t+1} | s_t = s] = -5$$

This  $R_s$  is simply the Reward after being in  
state  $s$ , because rewards are assigned / produced  
after the state at time  $t+1$ .

$\rightarrow \gamma$  is the discounting factor



f) return

Since Rewards at a particular state are not simply what we wish to calculate but rather cumulative reward across the whole action chain

$y_t$  is the total discounted reward from time-step  $t$

$$y_t = R_{t+1} + \gamma R_{t+2} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

The discounting factor we choose is  $\gamma \in [0, 1]$

if  $\gamma=0$  then implies we care about the immediate reward

if  $\gamma=1$ , that's maximally fanatical & it means we care about all rewards going into the future equally, assuming there is a terminal state in our markov process of states

The value of receiving a reward  $R$  after  $k+1$  time steps is  $\gamma^k R$

g) Value function

The value function  $v(s)$  gives the long term value of state  $s$

Formally it is simply the expected reward of starting in a state  $s$

so if we are dropped in a particular state what is the reward will be got from being in that state to the end till the process terminate

The state value function  $V(s)$  of an MRP is the expected return starting from state  $s$

$$V(s) = E[R_t | S_t = s]$$

The expectation exists because the environment is stochastic, so the evolution of states goes one way one step or another. So for one particular state we want to calculate the value over as as expectation i.e. as an average over all the possible random variable of all possible states, when the starting state is  $s$

### Example

Sample Returns for Our Student MRP

Starting from  $S_1 = C_1$  with  $\gamma = \frac{1}{2}$

For Sample  $\Rightarrow$

$C_1, C_2, C_3$  Pass Sleep

$$y_t = R_{t+1} + \gamma R_{t+2} - \dots$$

$$= R_{C1} + \gamma R_{C2} + \gamma^2 R_{C3} + \gamma^4 R_{\text{pass}} + \gamma^5 R_{\text{sleep}}$$

$$= -2 + \frac{1}{2} \times (-2) + \frac{1}{4} \times (-2) + \frac{1}{8} \times 10 + \frac{1}{16} \times 10$$

$$= -2.25$$

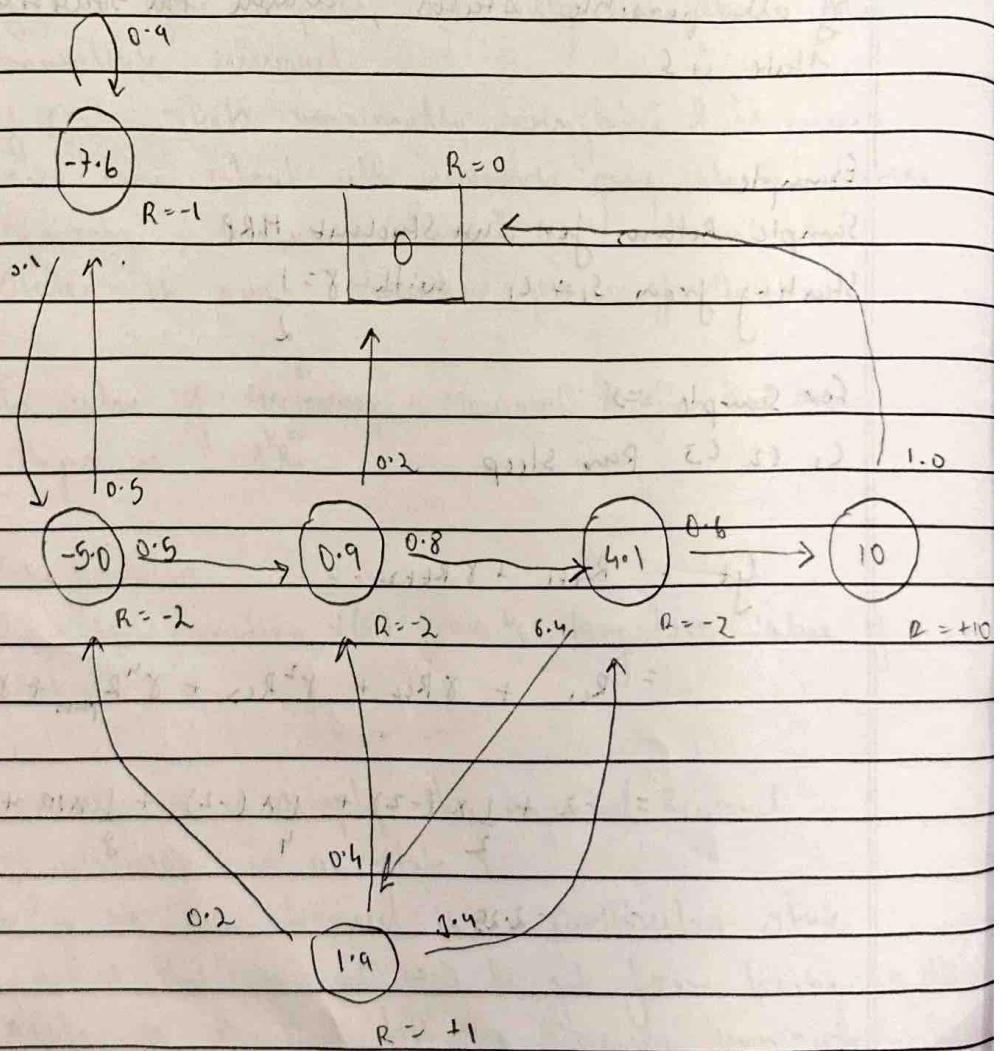
So, now one way to get the value of being in this state ( $C_1$ ) is to take such samples like:

$(1, FRFB, 1, C_2, SLEEP)$  or

$(1, C_2, 13, PVB, 12, C_3, Pm, SLEEP)$  or etc..

Then after calculating  $G_t$  we can (for each sample process) average them, to get an estimate value of the state  $C_1$ .

Example



Now to calculate the value function of a particular state (sum up in each state the numbers which the shape is the calculated value function for that state)

So for running (lens C<sub>2</sub> state)

$$\begin{aligned}
 \text{The } V(C_2) &= R_{C_2} + \gamma (0.8 \times V(C_3) + 0.2 (V(\text{sleep}))) \\
 &= -2 + 0.9 (0.8 \times 4.1 + 0) \\
 &= 0.952 \\
 &\approx 0.9
 \end{aligned}$$

To calculate  $V(C_2)$  we would need  $V(C_3)$  & hence each value of the states that  $C_2$  depends on in the future.

This formalizing can be better expressed & illustrated by the Bellman equation explanation in the following parts.

#### ii) Bellman Equation

The value function can be decomposed into two parts

→ The immediate reward  $R_{t+1}$

→ discounted value of a successor state  $\gamma, V(S_{t+1})$

$$\begin{aligned}
 V(s) &= E[Y_t | S_t = s] \\
 &= E[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | S_t = s] \\
 &= E[R_{t+1} + \gamma (R_{t+2} + \gamma R_{t+3} + \dots) | S_t = s] \\
 &= E[R_{t+1} + \gamma V(S_{t+1}) | S_t = s]
 \end{aligned}$$

So now, in the previous example

calculating the value of state  $C_1$  i.e.  $V(C_1)$

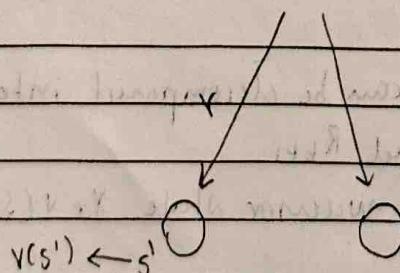
is simply the reward of being in that state which is  $-2$  ( $R_{t+1}$ )  $\pm$  then

discounting factor ( $0.9$ ) times the value of the next state, here since there are possible two next states thus the model term is  $\delta [0.8 \times V(C_2) + 0.2 \times V(Sleep)]$

### 1) Understanding Bellman Equations.

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

$$V(s) \leftarrow s \circ$$



The circles represent the states. At each state we have value functions for that state.

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

- 5) The Bellman equation can be expressed similarly using matrices as well;

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

where  $V$  is a column vector with one entry per state

$$\begin{bmatrix} V_1 \\ \vdots \\ V_n \end{bmatrix} = \begin{bmatrix} R_1 \\ \vdots \\ R_n \end{bmatrix} + \gamma \begin{bmatrix} P_{11} & \dots & P_{1n} \\ \vdots & \ddots & \vdots \\ P_{n1} & \dots & P_{nn} \end{bmatrix} \begin{bmatrix} V_1 \\ \vdots \\ V_n \end{bmatrix}$$

- K) Solving the Bellman equation

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

$$(1 - \gamma P)V = R$$

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

The computational complexity for  $n$  states is  $O(n^3)$

## L) Markov Decision Process

An MDP is a Markov Reward Process with decisions. It is an environment in which all states are Markov.

A markov decision process is a tuple  $\langle S, A, P, R, \gamma \rangle$

→  $S$  is a set of all possible states.

→  $A$  is our finite set of actions.

It is a discrete set of actions of finite size

→  $P$  is a state transition probability matrix,

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

Our transition probability matrix now depends on what action we take

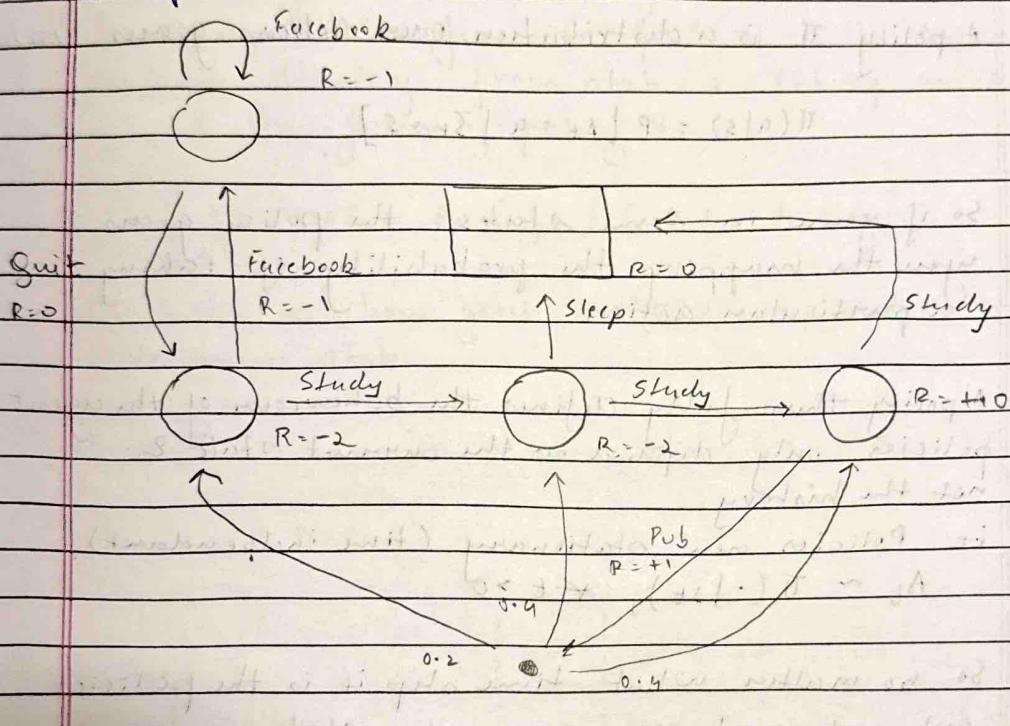
So now there will be a separate matrix

of  $P_{ss'}$  for each action  $a$

→  $R$  is a reward function,

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

→  $\gamma$  is the discounting factor

Example

Now we have more control over the MDP with actions through our path, also rewards are now attached to actions.

Our goal now is to find the best path through our decision making process that maximizes the sum of reward we get.

## M) Policies

A policy  $\Pi$  is a distribution over actions given state

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

So if you're in some states the policy gives you the mapping the probability of taking a particular action

A policy thus fully defines the behaviour of the agent policies only depend on the current state & not the history

i.e. Policies are stationary (time independent)

$$A_t \sim \Pi(\cdot | s_t), \forall t > 0$$

So no matter what time step it is the policies only depend on our given state

## N) Value Function

The state - value function  $V_\Pi(s)$  of an MDP is the expected return starting from state  $s$  & then following policy  $\Pi$

$$V_\Pi(s) = E[\gamma^t R_{t+1} | S_t = s]$$

So now  $V_\Pi(s)$  calculates the rewards dependent on a specific policy given we are in a particular state.

Now we sample all action according to the policy  $\Pi$

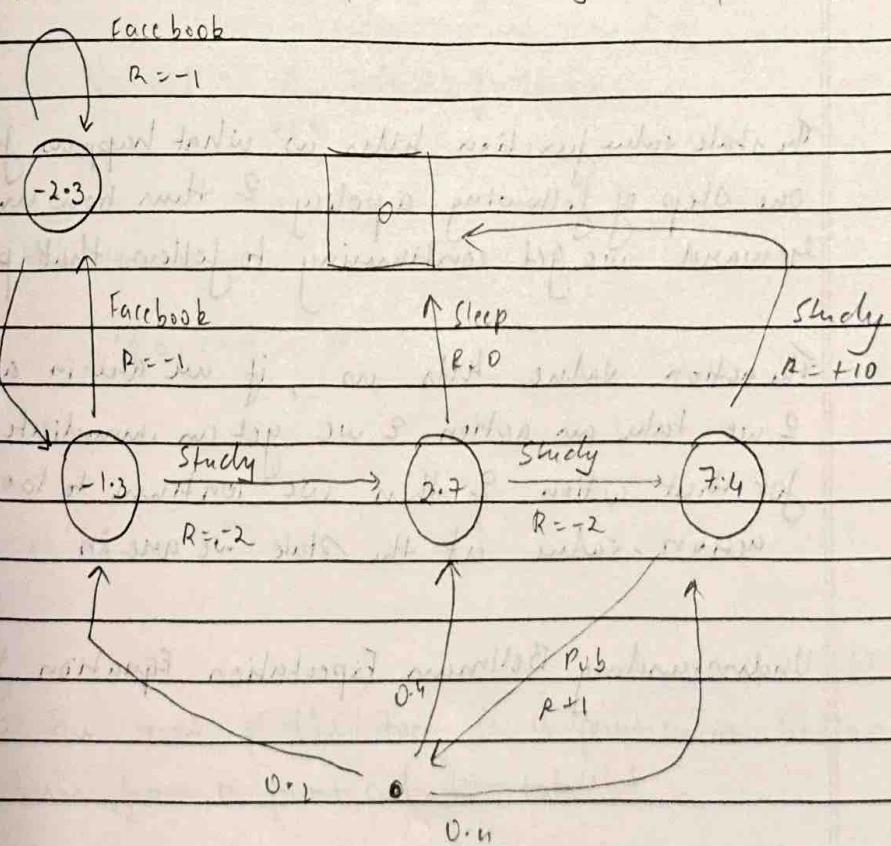
We also define another value function:

The Action-value function,  $q_{\pi}(s, a)$  is the expected return starting from state  $s$ , taking an action  $a$ , & then following policy  $\pi$ .

$$q_{\pi}(s, a) = E_{\pi} [g_t | s_t = s, a_t = a]$$

$q_t$  tells us how good is taking an action in a particular state.

Example:  $V_{\pi}(s)$  for  $\pi(a|s) = 0.5 ; \gamma = 1$



Here we are fixing our possible policy to be a so-so choice.

### o) Bellman Expectation Equation

The state value function can again be decomposed into immediate reward plus discounted value of a successor state,

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

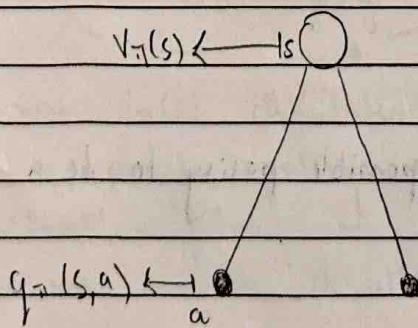
The action-value function can similarly be decomposed,

$$q(s, a) = \underset{\pi}{E} [R_{t+1} + \gamma q_{\pi}(s_{t+1}, a_{t+1}) | s_t = s, a_t = a]$$

The state value function tells us what happens for one step of following a policy & then how much more reward we get continuing to follow that policy

The action value tells us, if we are in a state  $s$  & we take an action  $a$  & we get an immediate reward for that action & then we continue to look at the action-value at the state we are in

### o) Understanding Bellman Expectation Equation for $V^*$

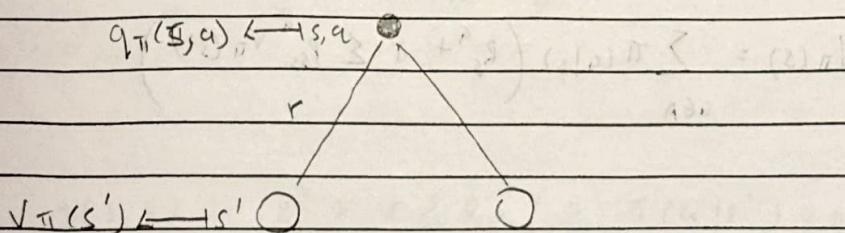


The white circle represents the state & the black dots represent the actions we can take from the given state

The probability of taking an action over a possible action is defined by our policy and for each of the action we take there is a  $q$  value associated to it & telling us how good it is taking action  $a$  from the given state we are in, & we simply average this across all actions

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

Now the concern,

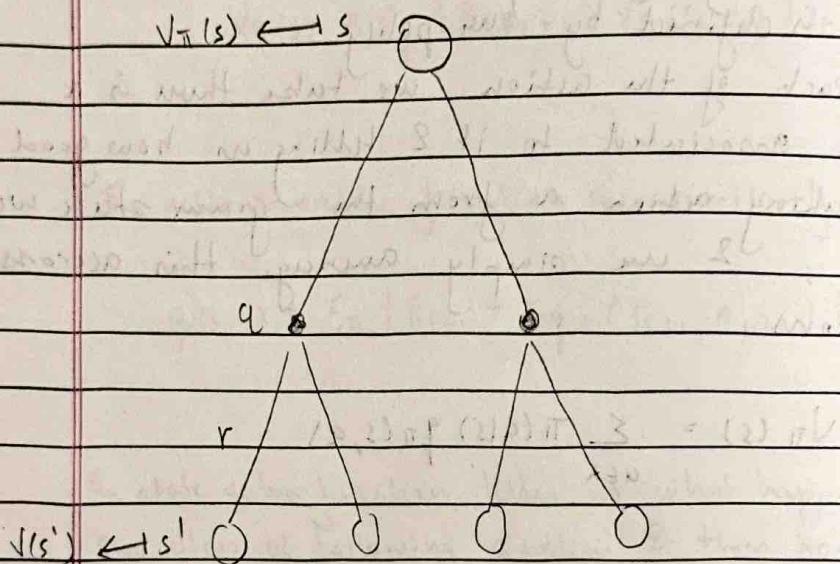


Now the root of the tree is a particular action we take from a particular state

Now for all the possible states we can go in taking a particular action, we average over the value (state value function) of being in the new state

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

Now to imagine how both work

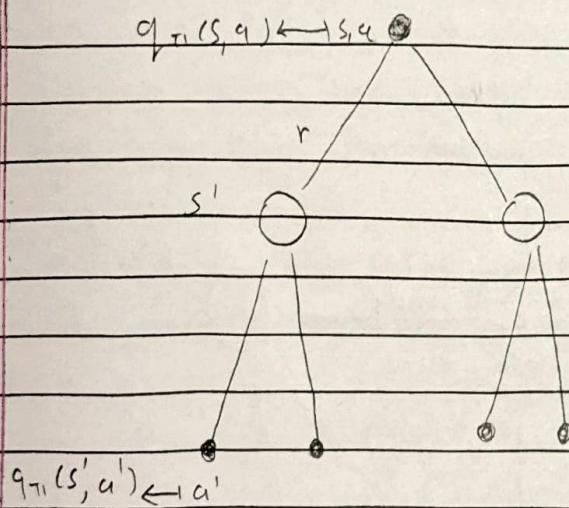


$$V_{\pi}(s) = \sum_{a \in A} \pi(a|s) \left( R_s^a + \gamma \sum_{s' \in S} P_{ss}^a V_{\pi}(s') \right)$$

So at the root of the tree we have a value function for a particular state, to get an idea of how good it is to be in that state we can average over first the action given our policy & then given that action the possible states we can reach & their value functions (state value functions). So this can be thought as it takes some

probability to take an action , & from that action to go to a state it takes/get a reward & then we also add a discounting factor to the state-value of the states we can reach from the particular action

To calculate action value



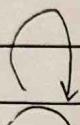
$$q_\pi(s, a) = R_s + \gamma \sum p_{ss'} \left( \sum_{a' \in A} \pi(a'|s') q_\pi(s', a') \right)$$

So this is simply for the value of taking a particular action from a state. we look at all possible states we can reach from the action in a given state, we get the  $R_s'$  since we take that action, within each state we discount by the discount factor the probability of going into

that state 2 then simple average over all possible actions out of that state given their state action values.

Q). Example

Facebook



$R = -1$

(-2.3)

0



Quit

$R = 0$

facebook

$R = -1$

Sleep

$R = 0$

(-1.3)

Study

$R = -2$

Study

$R = -2$

Study

$R = +10$

$P_{ab} = 1, R = +1$

0.4

0.2

0.4

Now let's calculate & verify our value function for the 7.4 state.

$$\text{Value(State)} = \frac{1}{2} \text{Study} + \frac{1}{2} \text{Pub}$$

$$\begin{aligned}
 &= \frac{1}{2} \times ( +10 + 0.8 \times 1 \times V_{\pi}(Sleep) ) \\
 &\quad + \frac{1}{2} \times ( +1 + 0.8 (0.6 \times 7.4 + 0.4 \times 2.7 + 0.2 \times (-1.3)) ) \\
 &= 5 + 2.34 \\
 &\approx 7.34 \\
 &\approx 7.4
 \end{aligned}$$

### A) Optimal Value Function

There could be multiple policies we follow in our markov chain, we wish to focus on the one that maximizes the reward.

The optimal state-value function  $V^*(s)$  is the maximum value function over all policies

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

$V^*(s)$  tells us what the maximum possible reward we can extract from the system.

The optimal action-value function  $Q^*(s, a)$  is the maximum action-value function over all policies

$$Q^*(s, a) = \max_{\pi} Q_{\pi}(s, a)$$

$Q^*(s, a)$  tells us the max reward we can extract starting in state  $s$  & taking action  $a$ , so given that we commit to an action which the most possible reward we can get from that

point onwards.

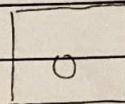
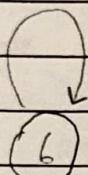
Example:

$$q^*(s, a) \text{ for } r=1$$

Facebook

$$R = -1$$

$$q^*_s = 5$$



(Quit

$$R = 0$$

$$q^*_s = b$$

Facebook

$$R = -1$$

$$q^*_s = 5$$



Study

$$R = -2$$

$$0.2$$

Sleep

$$R = 0$$

$$q^*_s = 0$$

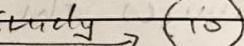
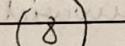
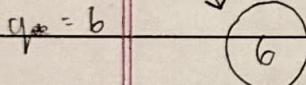
$$0.4$$

$$P_{ab}$$

$$R = +1$$

$$G_a = q^*_s \cdot P_{ab}$$

$$0.4$$



$$0.2$$

$$0.4$$

$$0.2$$

$$0.4$$

Each state denoted by a circle contains the optimal state-value, which tells us the maximum possible value of being in that state.

The  $q^*(s, a)$  values tell us that from each state given an action, which one brings the most value.

## 5) Optimal Policy

We can define a partial ordering over policies

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

## Theorem

for any Markov Decision Process

- There exists an optimal policy  $\pi^*$  that is better than or equal to all other policies,  $\pi^* \geq \pi$ ,  $\forall \pi$
- All optimal policies achieve the optimal value function,  
 $V_{\pi^*}(s) = V^*(s)$
- All optimal policies achieve the optimal action value function,  
 $q_{\pi^*}(s, a) = q^*(s, a)$

## 7) Finding an Optimal policy

An optimal policy can be found by maximizing over  $q_{\pi^*}(s, a)$ ,

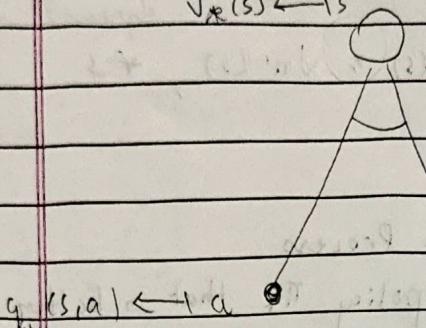
$$\pi^*(a|s) = \begin{cases} 1 & \text{if } a = \underset{a \in A}{\operatorname{argmax}} q^*(s, a) \\ 0 & \text{Otherwise} \end{cases}$$

So if we are in state  $s$  we pick the action with the max  $q^*$  value with a full probability

There is always a deterministic optimal policy & once we have  $q^*(s, a)$  we have the optimal policy

### v) Bellman Optimality Equation

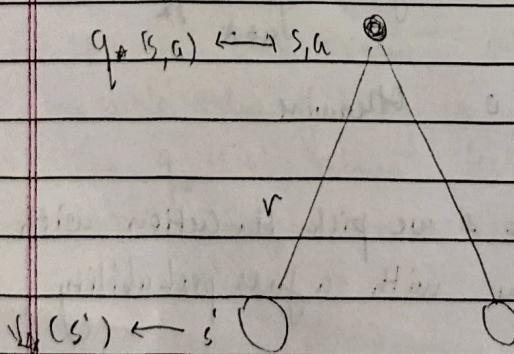
$$V_{*}(s) \leftarrow 1 s$$



Now for a given state  $s$  to get the  $V_*(s)$   
 we look at the possible actions to be taken  
 from each the given state , & we look at  
 the  $q_*(s,a)$  for each of them & then instead  
 of averaging over all of them we take the max

$$V_*(s) = \max_a q_*(s,a)$$

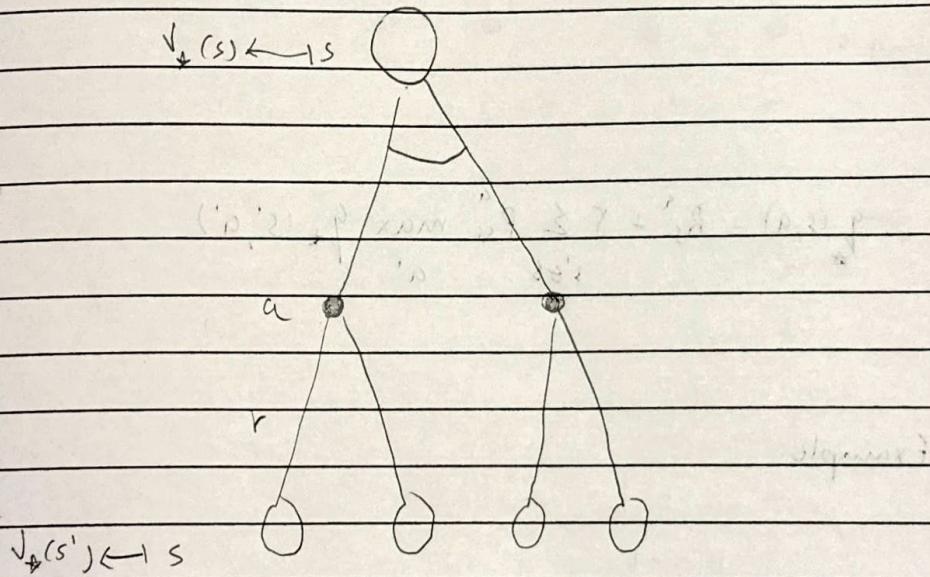
Now for the action  $a$ ,



After taking an action, we could end up in a possibility of states, & wif we arrange over the value function of each state that tells us how good our action is.

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

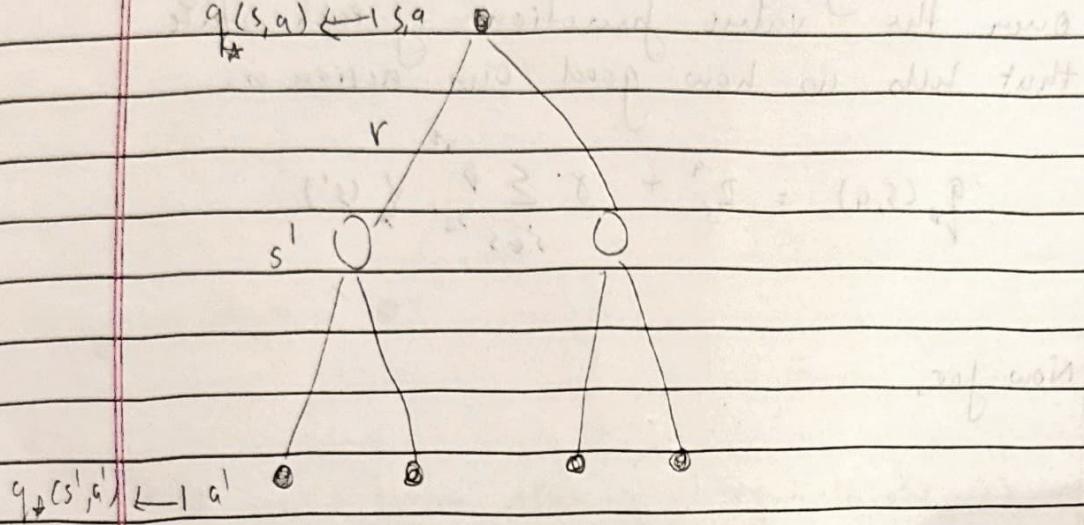
Now for,



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

For finding the optimal action-value

$$q(s, a) \leftarrow 1 s, a$$



$$q_{\pi}(s', a') \leftarrow 1 a'$$

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

Example

