

1 MARKOV DECISION PROCESSES

1

1.1 Introduction

- In reinforcement learning an AI learns how to optimally interact in a real-time environment using time delayed labels called rewards as a signal
- The markov decision process is a mathematical framework for R.L using states, and actions. and rewards
 - ↳ Through interacting w/ the environment, an AI will learn a policy which will return an action for a given state w/ highest reward.

1.2 The Bellman Equation

- Important Bellman Concepts

- STATE: numeric rep of what the agent is observing @ particular point of time in environment.
Eg: Raw pixels on a screen (in a game environ)
- ACTION: Input, agent provides to environ: calc by applying policy to current state.
Eg: Combination of control buttons being pressed @ given time.
: Analog i/p from joy stick
- REWARD: off Feed back signal : reflect how well the agent is performing.
Eg: coins collected, enemies killed ...
- In R.L, given current state we are, choose optimal action which will max the long-term expected reward provided by environ.
- Dr. Richard Bellman: Father of dynamic programming - class of algos simplify complex prob by breaking into sub prob and solving,

- What do Bellman Eq answer?

- Given state one is in, assuming best possible action taken @ each step what long term reward can one expect?
- What is the VALUE of the STATE?
- ∴ The bellman eq helps eval expected reward relative to adv or disadvantage of each state.

- The equation

$$V(s) = \max_a (R(s,a) + \gamma V(s'))$$

$V(s)$: Value of a given state

$R(s,a)$: Reward of optimal action a in state s

γ : Discount factor (diminishes reward over time)

$V(s')$: Value @ the next state.

Value of given state = max action (of all actions avail in the state, pick action that max value) · Reward of optimal action in state s to which we add so that the value of the next state times γ . It's where the recursive behaviors

comes about, this continues until terminal state is arrived

- The above equation would entail us exploring every possible action to arrive @ best long term reward. This type of bush force techniques might be obsolete w/ tasks that have too many moves to account for.

• Example :

→ if $\gamma = 1$ the $V=1$ @ all steps and agent will move around aimlessly.

→ γ : hyperparam

→ successful values b/w 0.9 & 0.99

→ small γ short-term thinking

→ large γ long-term thinking

$V=0.81$ $0+0.9(0.9)$	$V=0.9$ $0+0.9(0.9)$	$V=1$ $1+0.9(0)$	+1
$V=0.73$ $0+0.9(0.81)$		$V=0.9$	-1
$V=0.66$ $0+0.9(0.73)$	$V=0.73$	$V=0.81$	$V=0.73$

e.g. $\gamma = 0.9$

$$V(s) = \max_a (R(s,a) + \gamma V(s'))$$

1.3. Reading Assignment

3

- The rubic's cube

- State $s = \{r, b, g, 0, 0, y\}^{6 \times (3 \times 3)}$
- $X(s)$: avail actions to agent @ current state s
- Time t assumed to be discrete. $\in \{1, 2, \dots\}$
- Seq of steps $S_1 A_1 S_2 A_2 \dots$
 \hookrightarrow The trajectory is described what's happening but no info on the goal of agent.

- Feedback should be intro. to do help agent understand connection b/w its current actions & achievement of goal. In RL this takes the form of a scalar signal called Reward.

\hookrightarrow Seq: $S_1 A_1 R_2 S_2 A_2 R_3 \dots$

\rightarrow Expected tot reward ($+1$ state following action solved cube -1 otherwise)

$$\sum_{T=1}^T R_{t+1} = -(T-1) + 1 = -T + 2$$

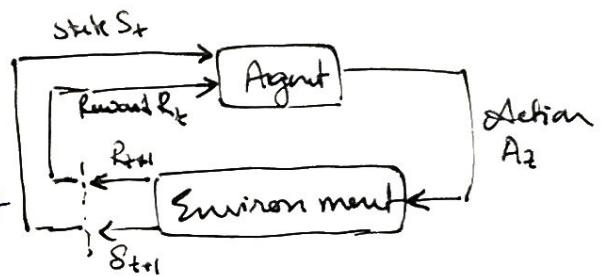
T : no. of moves agent try to reach solution

Lower T higher the reward: push the agent to do better.

- The common principle underlying R-L techniques is the reward hypothesis
 \hookrightarrow goals and purposes can be well thought as maximization of expected value of cumulative sum of rewards

- Markov Decision Process (MDP)

- Most of systems messier than rubic's cube: evolution of a system might be stochastic instead of deterministic
 \hookrightarrow A specific pair of (state,action) may lead to diff rewards ^{and give general formula} each w/ ^{its} own prob. We need to take all possibilities into consideration



- R.L ^{systems} is approached as a stochastic process.
- ↳ Deterministic models: S/P of model fully determined by param values and initial condi
- ↳ Stochastic models: process inherent randomness: same set of params and initial condi leads to diff results
- Seq of action states {rewards}: $S_1 A_1 R_1 S_2 \dots$ \leftarrow Random variables
- The environment fully known once properly specified.

↳ prob of receiving reward r @ time $t+1$ knowing system history upto that moment:

$$P(R_{t+1} = r | A_t = a_t, S_t = s_t, A_{t-1} = a_{t-1}, \dots, S_1 = s_1)$$

for each $t \in \mathbb{N}$, for $s_t, s_{t+1}, \dots, s_i \in S$ & for all $a_i \in A(s_i)$
 $i \in \{1, \dots, t\}$

↳ Prob of ending up in state s @ time $t+1$ knowing system history upto that moment:

$$P(S_{t+1} = s | A_t = a_t, \dots, S_1 = s_1)$$

↑ for ...

↳ We will always assume rewards are markovian: only dependent of on latest pair of state-action.

$$P(R_{t+1} = r | A_t, \dots, S_1 = s_1) = P(R_{t+1} = r | A_t = a_t, S_t = s_t)$$

often random variable S_t is also assumed to be markovian

↳ The environment is responsible for behaviour of states } rewards.

- The Agent is fully determined by prob of action a taken at time t ; knowing system history upto that moment.

$$P(A_t = a | S_t, A_{t-1}, \dots, S_1)$$

↳ Complete specification of agent behaviour is called Policy.

- Policy

• Notations

→ history of system; seq of states & actions observed so far.

$$h_t := (s_1, a_1, s_2, a_2, \dots, a_{t-1}, s_t)$$

$$s_i \in S \quad a_i \in A(s_i) \quad \text{for all } i \in \{1, \dots, t\}$$

H_t → history as random variable

→ choosing action space doesn't depend on current state of the system. (make sense, actions are always taken with intent to max reward @ next state)

$$H_1 := S$$

$$H_t := S \times A \times \dots \times S = S \times \prod_{i=1}^{t-1} A \times S \quad \text{if } t \in \{2, 3, \dots\}$$

$$H_1 := S$$

$$H_t := H_{t-1} \times A \times S \quad \text{if } t \in \{2, 3, \dots\}$$

Recursion

$$\therefore \text{history of system upto time } t: h_t = (s_{t+1}, a_{t+1}, s_{t+2})$$

• Decision Rules $\{d_t\}$ is the decision rule as it models how agent reaches acc

We know $P(A_t = a_t | S_t = s_t, A_{t-1} = a_{t-1}, \dots, S_1 = s_1)$ to current simulation experience?

→ Using new history information $P(A_t = a_t | H_t = h_t)$
which is prob of agent performing action @ time t ↑
knowing what has happened so far in terms of states & actions.

→ formalize it as function d_t from H_t to $P(A)$, The
Space of prob distribution over A :

$$d_t: H_t \rightarrow P(A)$$

→ We shall denote by $b_{d_t}(a | h_t)$ or $b_{d_t(h_t)}(a)$ The
prob of taking action a if our system history upto
time $t = h_t$ ie

$$P(A_t = a_t | H_t = h_t) = b_{d_t(h_t)}(a_t)$$

• Policy

→ It is the recipe underlying the behaviour of our agent: what it's programmed to do under certain circumstances.

→ Decision rules provide a quick way to formally define it.
 ↳ policy π : seq of decision rules.

$$\pi = (d_1, d_2, \dots) \quad \text{w} \quad d_i : H_i \rightarrow P(A) \text{ for each } i \in \{1, 2, \dots\}$$

→ Reward hypothesis: The goal is the maximization of expected value of the cumulative sum of a received scalar signal.

• Returns

→ In RL every action is followed by reward. (+ve / -ve) → measures overall performance of a policy π of rewards \Rightarrow called Return

→ If R_{t+1} reward after action A_t ; cumulative return at t , G_t

$$G_t := \sum_{k=t+1}^T R_k \quad T: \text{random variable denoting end of time of our simulation}$$

→ Sometimes it's the case where $T = +\infty$: exploration game.

↳ The issue can be solved by intro discount rate $\gamma \in [0, 1)$

$$\text{Discounted Cumulative Return} \quad G_t := \sum_{k=t+1}^T \gamma^{k-t-1} R_k$$

$\gamma \in [0, 1)$ allows \mathbb{E} w.r.t. to closer rewards and ensures overall reward is finite if rewards are bounded.

$$|G_t| \leq \sum_{k=t+1}^{+\infty} |\gamma^{k-t-1} R_k| = \sum_{k=t+1}^{+\infty} \gamma^{k-t-1} |R_k| \leq M \sum_{k=t+1}^{\infty} \gamma^{k-t-1} = \frac{M}{1-\gamma}$$

assuming there exists $M > 0$ such that $|R_k| \leq M \forall k$

• State-value function

→ let $\pi = (d_1, d_2, \dots)$ be policy followed by agent

↳ G_t = discounted cumulative return from time t .

↳ S_0 = initial state of environment. Then expected value of cumulative sum of received scalar signal \mathbb{E}_{π}^S $\mathbb{E}_{\pi}^S(S_0) := \mathbb{E}_{\pi}[G_t | S_0 = s]$
 also called S.v.f. of policy π

- Expected Reward

→ Expected value of 1st reward, R_1 under policy π

$$E_{\pi}[R_1 | S_1 = s_1]$$

↳ Expectation of $R \cdot V$ = sum of all possible values wtd by its probability

$$= \sum_{r \in R} r P(R_1 = r | S_1 = s_1)$$

R is a set of possible rewards.

↳ Using law of total prob → this can be decomposed w.r.t set of possible actions available to agent (A_1)

$$= \sum_{a \in A} P(R_1 = r | A_1 = a, S_1 = s_1) P(A_1 = a | S_1 = s_1)$$

↳ But $P(A_1 = a | S_1 = s_1)$ is fully determined by agent's 1st decision rule $d_1 : b_{d_1(s_1)}(a)$

$$\therefore E_{\pi}[R_1 | S_1 = s_1] = \sum_{a \in A} b_{d_1(s_1)}(a) \left(\sum_{r \in R} r P(R_1 = r | A_1 = a, S_1 = s_1) \right)$$

$$= \sum_{a \in A} b_{d_1(s_1)}(a) \cdot \underbrace{E(R_1 | A_1 = a, S_1 = s_1)}_{\text{agent}} \quad \underbrace{- (1)}_{\text{environment}}$$

The policy determines the prob to choose an action, while the environment is responsible for the expected return awaiting the agent if it were to choose a specific course of action.

- Expected Return : to derive formula for expected discounted return.

→ All expectation is linear ie $E[\alpha A + \beta B] = \alpha E(A) + \beta E(B)$

A & B are random variables, α & β real nos.

→ Then

$$\begin{aligned} v_{\pi}'(s_1) &= E_{\pi}[G_1 | S_1 = s_1] \\ &= E_{\pi}[R_1 + \gamma G_2 | S_1 = s_1] \\ &= E_{\pi}[R_1 | S_1 = s_1] + \gamma \underbrace{E_{\pi}[G_2 | S_1 = s_1]}_{\text{already derived formula.}} \quad - (2) \end{aligned}$$

→ Using law of tot prob and def of expected value

$$\square = \sum_g \sum_{a \in A} P(G_2 = g | A_1 = a, S_1 = s_1) b_{d, C(g)}(a)$$

→ Repeating this again: this time w.r.t S_2 :

$$P(G_L = g | A_1 = a, S_1 = s_1) = \sum_{\delta \in S} P(G_L = g | S_L = \delta, A_1 = a, S_1 = s_1) \cdot P(S_L = \delta | A_1 = a, S_1 = s_1)$$

$$= \sum_{\delta \in S} P(G_L = g | H_2 = (S_1, a, \delta)) \cdot P(H_2 = (S_1, a, \delta))$$

→ Putting everything together, we get:

$$\begin{aligned}
 \square &= \sum_g g \sum_{a \in A} \sum_{s \in S} P(G_1 = g | H_2 = (s_1, a_1, s)) \cdot b_{d_1(s_1)}^{(a)} \cdot P^{(u)} \\
 &= \sum_{a \in A} \sum_{s \in S} \left[\sum_g g P(G_1 = g | H_2 = (s_1, a_1, s)) \right] \cdot b_{d_1(s_1)}^{(a)} P^{(u)} \\
 &= \sum_{a \in A} \sum_{s \in S} E_{\pi} [G_1 | H_2 = (s_1, a_1, s)] \cdot b_{d_1(s_1)}^{(a)} \cdot P^{(u)}
 \end{aligned}$$

The diff b/w this and S.R.F is that we are computing the expected discounted cumulative return starting from $t=2$ instead $t=1$ as well as conditioning on H_2 instead of $S_1 = H_1$.

→ formalizing this by introducing the history-value function
for policy π @ time t .

$$V_{\pi}^+(h_t) := E_{\pi}[G_{t+1} | H_t = h_t]$$

↳ Plugging this into equation □

$$\boxed{Q} = \sum_{q \in A} \sum_{S \in S} \gamma_1^2((s_1, q, s)) \cdot b_{d_1(s_1)}^{(q)} \underbrace{\cdot P(\cdot)}_{\text{environment}} - \boxed{3}$$

from eq (1), (2) & (3) the general formula for ^{again} expected value - cumulative discounted return under policy π :

$$V_{\pi}^1(S_1) = \sum_{a \in A} b_{d,(S_1)}^{(a)} \left[ELR_2 | A_1=a, S_1=s_1 \right] +$$

$$\gamma \sum_{S_2 \in S} V_{\pi}^L((S_1, a, S)) \cdot P(S_2=s | A_1=a, S_1=s_1)$$

This is a generalized version of Bellman's Eq in the framework of MDP

Bellman Equation

→ Using techniques from before one can prove the following equation holds for $t \in \{1, 2, \dots\}$

$$v_t^*(h_t) = \sum_{a \in A} b_{d_t(h_t)} (a) [E[R_{t+1} | A_t = a, H_t = h_t] + \gamma \sum_{s \in S} v_{t+1}^*((h_{t+1}, s)) \cdot P(S_{t+1} = s | A_t = a, H_t = h_t)]$$

The problem here is: to compute v_t^* we need to know v_{t+1}^*
which in turn requires $v_{t+2}^* \dots$ no end assuming $T = +\infty$

↳ Things would simplify if we assume that the environment and the agent are markovian & stationary
ie $S_t \& R_t$ are only dependent on $A_{t-1} \& S_{t-1}$
while agent bases its A_t choice on the value of S_t
w/o taking into account current time instant t .
ie $\pi = (\pi_1, \pi_2, \dots)$

∴ Eqn (4) becomes

$$v_t^*(s_t) = \sum_{a \in A} b_{d_t(s_t)} (a) [E[R_t | A_t = a, S_t = s_t] + \gamma \sum_{s \in S} v_{t+1}^*(s) P(S_{t+1} = s | A_t = a, S_t = s_t)]$$

A policy value function is in fact the expected cumulative discounted reward following that strategy: if $v_t^*(s) \geq v_{t'}^*(s)$ for every $s \in S$ we can generally conclude π is a better policy than π' (the envr by returning actions to take for each state),

π^* is called the optimal policy that maximizes expected reward. Among all policies taken, the optimal policy is the one that maximizes the amt. of reward received or expected to receive. For an MDP there is no end of lifetime, hence you have to decide the time.

This policy is a guide telling which action to take for a given state. It's not a plan but returns the underlying plans of

1.4 MDP Decision Processes

10

- MDP is formally used to describe an environment for R.L where the environment is fully observable.

- The MDP

- The markov property states that: future is independent of past once given the present.

↳ a state S_t has Markov property iff

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

The state captures all relevant info from history

↳ For markov state S and successor state S' , the transition prob is defined by:

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

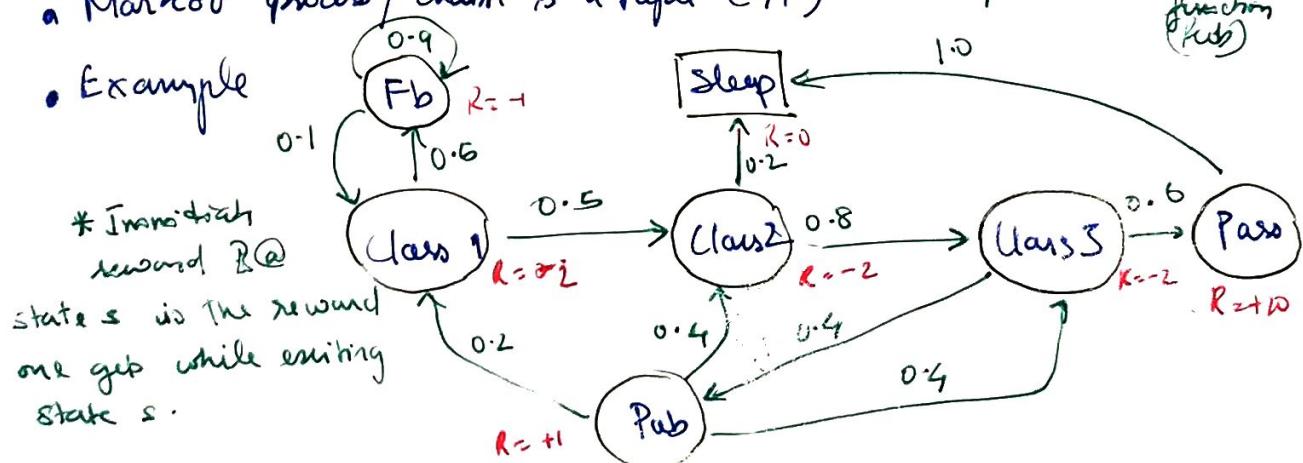
prob distribution over next possible successor state gives current state.

Matrix form of all possible transitions, where each row sums to 1

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

- Markov Process

- Memory-less, random process ie seq of random states S_1, S_2, \dots w/ markov property
- A Markov process / chain is a tuple (S, P) S : state space P : transition function $P(s, s')$
- Example



- Some example sample episodes 2 seq of states? or: 11
 - $C_1 \ C_2 \ C_3 \ Pass \ Sleep$
 - $C_1 \ C_2 \ C_3 \ Pub \ C_2 \ C_3 \ Pub \ C_2 \ Sleep$.
- The transition matrix for the same is

	C_1	C_2	C_3	Pass	Pub	FB	Sleep
C_1	0.5	0.8				0.5	
C_2		0.6	0.4			0.2	
C_3			0.6	0.4			
Pass	0.2	0.4	0.4				1.0
Pub	0.1					0.9	
FB							1.0
Sleep							

- Markov Reward Process

- The MRP is a tuple (S, P, R, γ)
 - S : a finite set of states
 - P : state transition prob matrix
 - $P_{SS} = P[S_{t+1} = S' | S_t = S]$
 - R : reward function, $R_S = E[R_{t+1} | S_t = S]$ } Reward for only 1 step.
 - γ : discount factor, $\gamma \in [0, 1]$
- Looking @ the MDP ex in page 10, rewards have arbitrarily been added to it: you don't like sitting in class $\therefore -2$, but if you pass all you get $+10$.
- We care about cumulative reward, as the goal is to max the return. \therefore The return G_t is the total discounted reward from time step t

$$G_t = R_{t+1} + \gamma R_{t+2} + \dots = \sum_{k=t+1}^{\infty} \gamma^{k-t-1} \cdot R_k$$
 - There is no expectation here cause G_t is just one sample from a set MRP of the rewards we get going through that seq.
- The value function is the long term value of being in a state. So were to drop you in state S , how much value will you get from there on. What's the tot reward you'll get from S onwards

$$V(S) = E[G_t | S_t = S] \quad (\text{using stochastic hence } E)$$

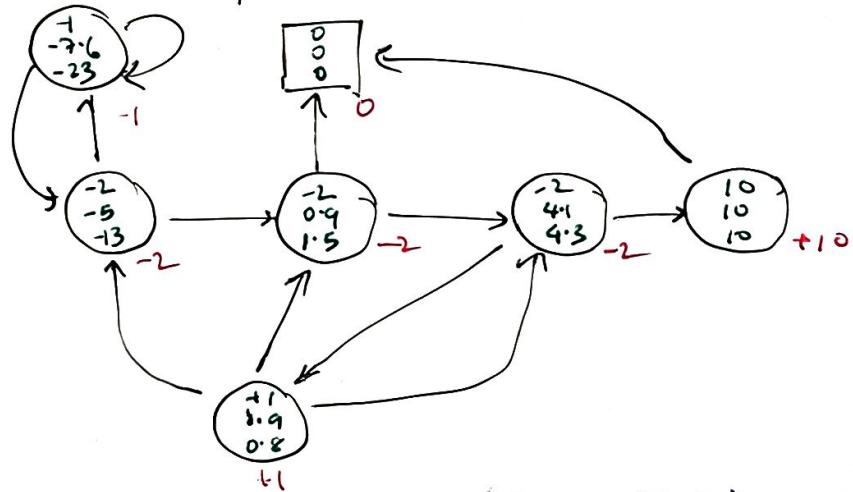
↳ Taking the student-class example: starting from $S_0 = C_1$ 12 and $\gamma = 0.5 \Rightarrow G_0 = R_2 + \gamma R_3 + \dots + \gamma^{T-2} R_T$

$$\begin{array}{l|l} \text{(C1 C2 C3 Pass Sleep} & V_1 = -2 - 2 \cdot 0.5 - 2 \cdot 0.25 + 10 \cdot 0.125 = -2.25 \\ \text{C1 FB FB C1 C2 Sleep} & V_1 = " - 1 \cdot 0.5 - 1 \cdot 0.25 - 2 \cdot 0.125 = -3.125 \\ \text{C1 C2 C3 Pass C2 C3 Pass Sleep} & V_1 = " \\ & = -3.41 \end{array}$$

We have 3 start off samples here, one way to estimate the value of the state would be take a kernel of samples and take the avg. and

∴ The returns of samples are random but the value function is not a random quantity, rather an expectation of these R.V

↳ The state-value w/ $\gamma = 0, 0.9, ? 1$



We observe as γ is increased (long sighted) we get better rewards given a particular state S_t as more wt given to the future rewards.

The Bellman Equation for MRPs

→ Value function decomposed into

↳ Immediate reward R_{t+1}

↳ Discounted value of successor state $\gamma V(S_{t+1})$

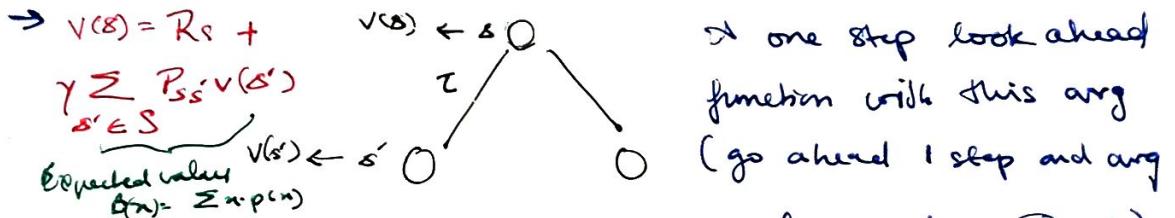
$$V(S) = E[G_t | S_t = s]$$

$$= E[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | S_t = s]$$

$$= E[R_{t+1} + \gamma (R_{t+2} + \gamma R_{t+3} + \dots) | S_t = s] \quad (\text{Immediate reward} + \text{discounted reward from next steps})$$

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

$$= E[R_{t+1} + \gamma V(S_{t+1}) | S_t = s] \quad \begin{aligned} &\text{The val. function now =} \\ &\text{expected immediate reward} \\ &\text{discounted value @ next state} \end{aligned}$$



all possible outcomes together which gives value function @ s)

In the future one will notice that building of algos will take into the idea of one step look ahead.

\hookrightarrow from the diagram in pag. 12: w/ $\gamma = 1$ & $C_3 = 4.3$

$$\begin{aligned} \text{using the above formula } (V_G) &= R_s + \gamma \sum P_{ss'} v(s') \\ &= -2 + 1(0.6 \cdot 10 + 0.4 \cdot 0.8) \\ &= 4.3 \end{aligned}$$

\rightarrow Matrix representation.

$$v = R + \gamma Pv$$

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

The value of state = immediate reward + the transition matrix + the value when we may end up

\rightarrow The Bellman Equation is a linear eq:

$$\hookrightarrow \text{can be solved directly } v = R + \gamma Pv$$

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

\hookrightarrow computational complexity = $O(n^2)$

\hookrightarrow Direct solution only possible for small MRP's

\hookrightarrow For large MRP's \rightarrow Dynamic Programming \rightarrow Temporal difference learning.
 \rightarrow Monte Carlo evaluation

- Markov Decision Process

Until now one scally hasn't had agency to take actions or make decisions. So let's introduce decisions.

∴ MDP is MRP with decisions.

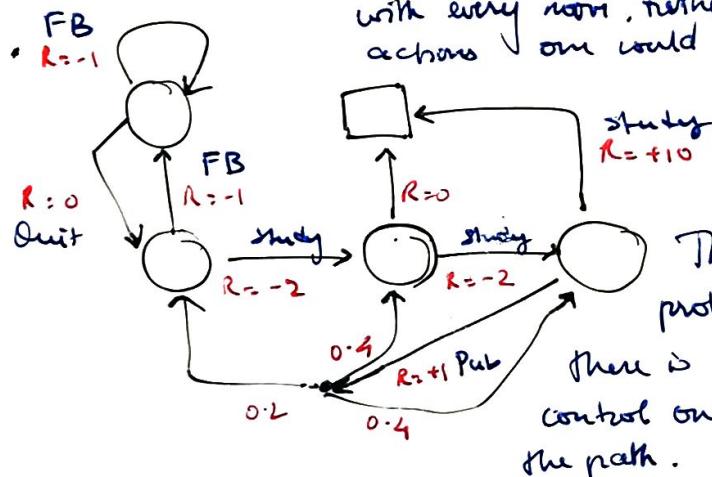
Hence MDP is a tuple (S, A, P, R, γ) where:

- S : finite set of states
- A : finite set of actions
- P : state transition prob form matrix

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

- R is a reward function $R_s^a = E[R_{t+1} | S_t = s, A_t = a]$
- γ : discount factor $\gamma \in [0, 1]$.

- Example: student MDP → So here there aren't probs attached with every move, rather there are actions one could decide on.



There are still some prob, but this time there is a lot more control one can exert on the path.

Policies

→ Distributions over actions given states

$$\pi(a|s) = P[A_t = a | S_t = s] \text{ ie if one is in}$$

some state s this distribution gives us the mapping -

→ A policy fully defines the behaviour of an agent

→ MDP policies depend on the current state (not history)

i.e. policies are stationary (time-dependent)

$$\pi(a|s), \forall t > 0$$

* There aren't any rewards in the equation because the state s fully characterizes future rewards. The markov property means → fully captures the evolution from state s onwards.

So we are looking for a policy, which given the state 15 one is in you'll want to take actions that'll get you the most future reward. (We don't care about rewards in the past as they are gone, all we care about is the rewards into the future)

→ Connection b/w MDP & MRPI we can always recover MRPI from MDP

↳ The state seq, $S_1, S_2 \dots$ is a markov process

↳ The state & reward seq, $S_1, R_1, S_2 \dots$ is a markov reward process $(S, P^\pi, R^\pi, \gamma)$

where

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

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

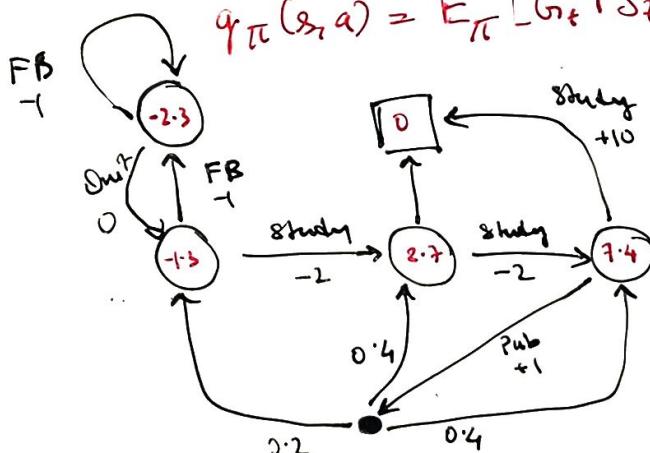
• Value Function

The state-value function $V_\pi(s)$ (value @ state s after having followed policy π) of an MDP is the expected return starting from state s and following policy π

$$V_\pi(s) = E_\pi [G_t | S_t = s] \quad \begin{matrix} \text{(interpretation when we} \\ \text{sample all policies)} \end{matrix}$$

→ Action value function $q_\pi(s, a)$ is the expected return starting from state s , taking action a , and then following policy π

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



$$V_\pi(s) \text{ for } \pi(a|s) = 0.5 \quad \gamma = 1$$

• Bellman Equation
Expectation

→ The state value function can again be decomposed into immediate reward + discounted reward

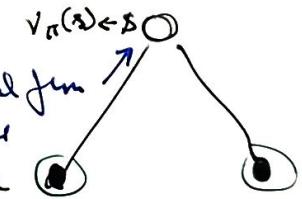
$$v_{\pi}(s) = E_{\pi}[R_{t+1} + \gamma v_{\pi}(s_{t+1}) | S_t = s]$$

→ Similarly, the action-value function can similarly be decomposed

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

↳ Tells us, if I'm in one state and an action is taken then some immediate reward is received for that action and then I look @ when I end up and ask what's the action value of the state I end up in under the action I would pick from that point on.

→ Understanding how v & q relate to each other. If we were in a state-val fun here. What it says is that we are going to avg over the $q_{\pi}(s, a) \leftarrow a$



state
action

actions we might take. There's some prob we might take this a action here, the prob of which is defined by the policy - (prob you'll go left, right...). For each of the actions we might take there's a q value telling how good is it to take $A_t = a$ from $S_t = s$. So what we are doing is a 1 step look ahead saying that the state value, how good is the state to look ahead 1 step, look at action values avg them together. and that tells us the value of being @ $s_t = s$

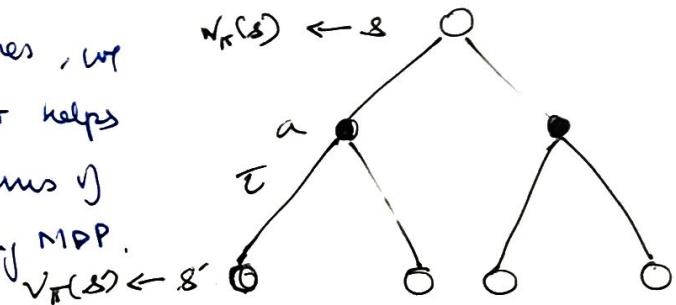
$$V(s) = \sum_{a \in A} \pi(a|s) q_{\pi}(s, a)$$

→ Understanding the opp step: starting w/ taking some action.

Now the root of this tree is $q_{\pi}(s, a) \leftarrow s, a$
 a state and we are considering
 a particular action we take from
 that state. If I'm in this state $v_{\pi}(s') \leftarrow s'$
 here, how good would it be to
 go right from that state? Then for v_{π} tells us how good is it
 to be in a particular state, q_{π} tells us how good is it to
 take an action from a given state.

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

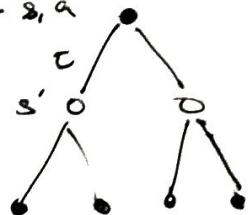
→ Putting stuff together, we
 get a recursion that helps
 us understand v in terms of
 itself → end up solving MDP.



At the roots of the tree, we have the val-function for
 a particular state - telling us how good is it to be in that
 state. The way we can understand that is by doing
 a 2 step look ahead - consider all the actions we might
 take next, consider all the things the environment might
 do to us, and for each of the things the environment would
 do there would be some success state we might end up in.
 And we'd want to know how good is it to be in that
 state and carry on with the policy. If we are carrying
 over policy, weight each of them $\pi(a|s)$ are by prob the
 policy will select, and among them $\pi(a|s)$ are key the
 transition prob that will end up getting blown in one
 direction or another. Avg all together and gives us the val
 of being at the root

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

→ One can do exactly the same thing $q_{\pi}(s, a) \leftarrow s, a$
 with action values (recursive relation)
 by sticking diagram other way
 around.

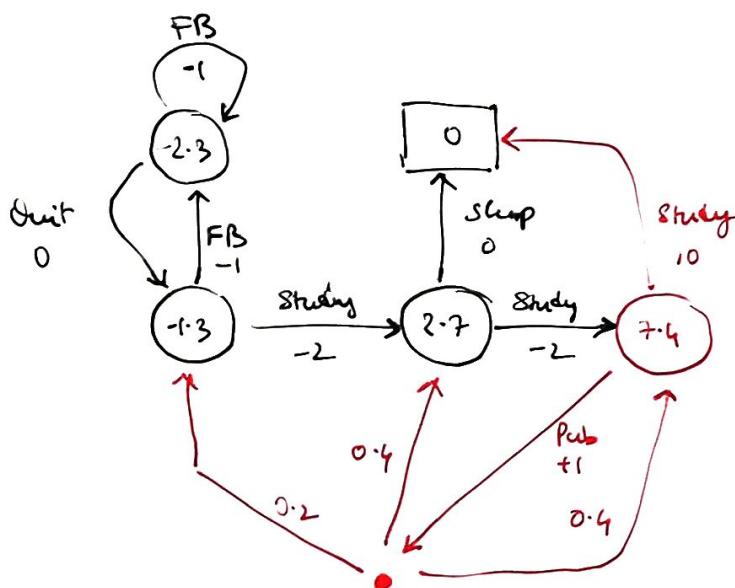


starting from a particular state or action we can look ahead 2 steps, consider which state we might end up and then consider from the state we end up which action might be taken next and then avg over it

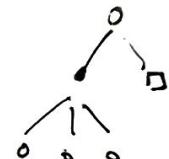
↳ In the previous diagram the state values relate to state values in next step and here the q values relate to the q values in the next step

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

→ Student MDP example



Considering only
 the red state;
 Class 3.
 $\pi(a|s) = 0.5$



$$\begin{aligned} V(s) &= \sum \pi(a|s) \left(R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a v_{\pi}(s') \right) \\ &= 0.5 \times 10 + 0.5 (1 + \gamma (0.2 \times -1.3 + 0.4 \times 2.7 + 0.4 \times 7.4)) \\ &= 7.4 \end{aligned}$$

Hence the bellman equation does really hold. This is all well and good but it doesn't tell us the best way to behave.

→ Matrix form of bellman eq

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

$$\text{or direct sol} \rightarrow V\pi = (1 - \gamma P^\pi)^{-1} R^\pi$$

- Optimal Value Function

The main problem would be finding the best behaviour using MDP.

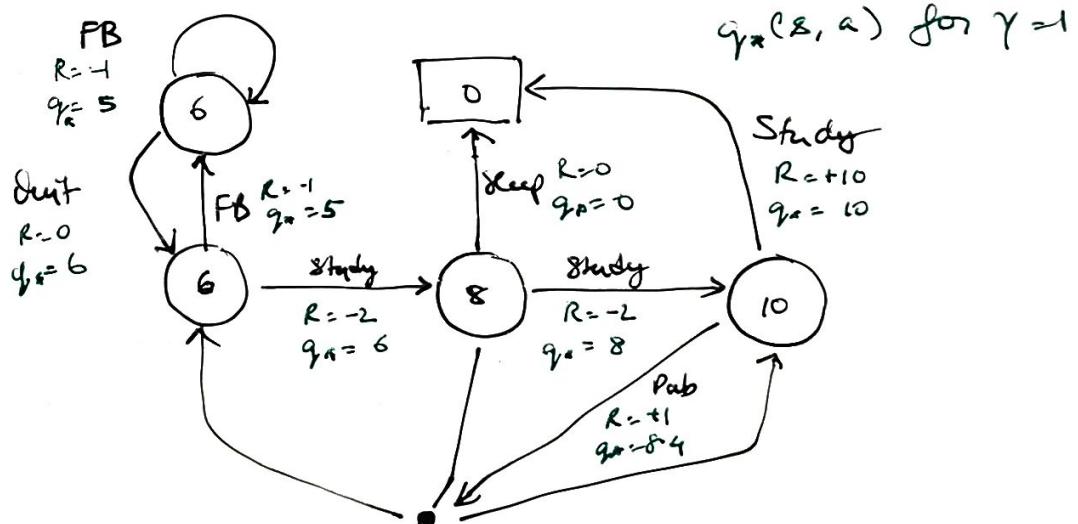
The optimal state-value function $V^*(s)$ is the max value function over all policies.

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

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

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

→ MDP is solved when we know the optimal value function $q^*(s, a)$ for $\gamma = 1$



→ Optimal policy: we want to know what is it for one policy to be better than another policy and define a notion of optimality.

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

Defining partial ordering of policies

↳ For any MDP :

(i) There exists optimal policy π^* that is better or equal to

all other policies $\pi \geq \pi^*, \forall \pi$

(ii) All optimal policies achieve the opt-val fun $V_{\pi^*}(s) = V^*(s)$

(iii) " the opt-act-val fun $q_{\pi^*}(s, a) = q^*(s, a)$

↳ Finding Opt policy : can be found by max over $q_\pi(s, a)$

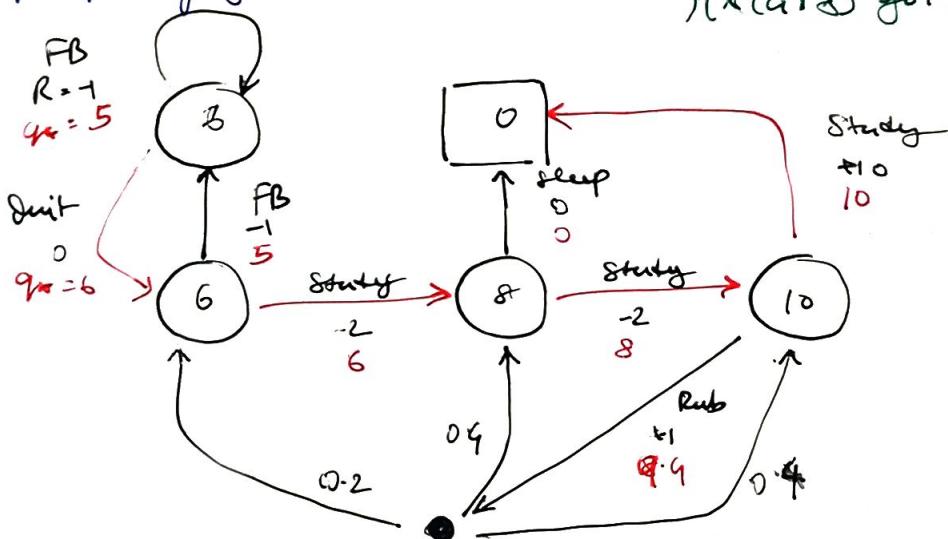
$$\pi^*(a|s) = \begin{cases} 1 & \text{if } a = \underset{a \in A}{\text{argmax}} q^*(s, a) \\ 0 & \text{otherwise} \end{cases}$$

* There is always a deterministic opt policy for any MDP

* If we know $q^*(s, a)$ we immediately have opt-policy

→ Opt-poly for student MDP

$\pi^*(a|s)$ for $\gamma = 1$



Red arc : optimal policy π^*

$$\text{Pub arc} : = 1 + \{0.2 \times 6 + 0.4 \times 8 + 0.4 \times 10\} \quad \left. \right\} \text{max} = 10$$

$$= 9.4$$

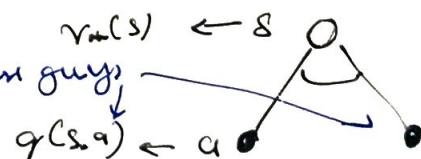
$$\text{Study arc} = 10 + 1(0) = 10$$

• Bellman Optimality Equation

→ Tells us how to solve MDP; how to relate opt-val fun to itself

→ Instead of taking Expval of these guys
we take the max

$$V_{\pi^*}(s) = \max_a q_{\pi^*}(s, a)$$



→ Optimality Eq for Q^* $q_{**}(s,a) \leftarrow s,a$

Now we want to know

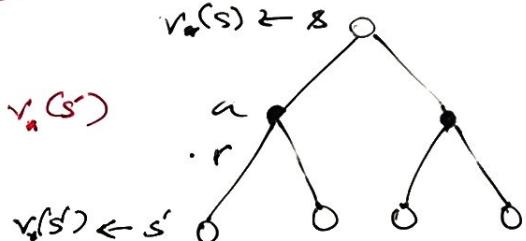
how good are the red ones $v_x(s') \leftarrow s'$

We don't control this part. But each of the state we may end up in would have one opt value

$$q_p(s,a) = R_s^a + \gamma \sum_{s' \in S} p_{ss'}^a v_x(s')$$

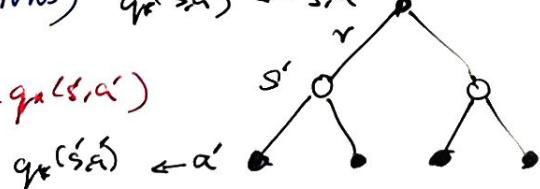
→ 2 step look ahead

$$v_x(s) = \max_a R_s^a + \gamma \sum_{s' \in S} p_{ss'}^a v_x(s')$$



→ 2 step look ahead (actions) $q_p(s,a) \leftarrow s,a$

$$q_p(s,a) = R_s^a + \gamma \sum_{s' \in S} p_{ss'}^a \max_a q_p(s',a')$$



→ Solving the Bellman Opt Eq

↳ The eq: is non-linear now

↳ No closed-form solution (in general)

↳ Many iterative sol: val iter, policy iter, Q-learning, SARSA

- Extensions to MDPs

- Infinite MDPs

→ countably infinite state and/or action spaces

- Partially observable MDPs

- undiscounted avg reward MDPs