

lec 2

MDP $\xrightarrow{\text{environment}}$
(FULLY OBSERVABLE)

Any partially observable problem can be completely converted to MDP

Recall

Markov Property where current state categorizes all relevant information

II State transition probability

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

now we can create a matrix from all s to all s'

$$P = \begin{bmatrix} P_{11} & \dots & P_{1n} \\ \vdots & \ddots & \vdots \\ P_{m1} & \dots & P_{mm} \end{bmatrix}$$

current next

II Markov Process / $\xrightarrow{\text{Chain}}$ sequence of states with
Markov Property

Type (S,P)

1 /

terminal state (self-loop) like a rest state
Sample episodes just mean sequence of states

VI Markov Reward Process (of a sample episode)
 chain with value judgements

$$(S, P, R, \gamma)$$

↑ { matrix } $\in [0, 1]$

↳ farsighted
 ↳ shortsighted / myopic

$$R_s = \mathbb{E}[R_{t+1} | S_t = s]$$

next time step pe you get reward of
 current step

VI Return G_t

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

the more close it is to 0: we prefer immediate rewards

" to 1: indifferent to when rewards arrive

γ - mathematically convenient

- avoid infinite returns in cyclic processes
- if financial, we may want immediate rewards
 coz interest

VI State-value Function: long term value of being in state s
 : expected return starting from state s

$$V(s) = \mathbb{E}_{\text{State}}[G_t | S_t = s]$$

Take one start state . for many episodes figure
value & average (one way)

Bellman Equation for MRP

$$v(s) = \mathbb{E} [G_t | S_t = s]$$

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

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

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

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

↑
immediate
reward

↑
value function of next state

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

$$\begin{bmatrix} v(c_1) \\ \vdots \\ v(c_n) \end{bmatrix} = \begin{bmatrix} R_1 \\ \vdots \\ R_n \end{bmatrix} + \gamma \begin{bmatrix} P_{11} & \dots & P_{1n} \\ \vdots & \ddots & \vdots \\ P_{n1} & \dots & P_{nn} \end{bmatrix} \begin{bmatrix} v(c_1) \\ \vdots \\ v(c_n) \end{bmatrix}$$

↑ reward from exiting state

$$v = R + \gamma P v$$

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

$$v = (1 - \gamma P)^{-1}R \longrightarrow O(n^3)$$

not practical

only works for small MDPs

for large MDP — DP, Monte Carlo, Temporal Difference

MDPs are MRP with decisions/actions

tuple (S, A, P, R, γ)

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

↑ ↗
finite set of states & actions state transition matrix
[it now depends on action so new matrix for each action]

$$\text{Policy } \pi(a|s) = P[A_t=a | S_t=s]$$

Markov property holds that only current policy matters

i.e. time independent i.e. same policy in each time step

New definitions

$$V_\pi(s) = \mathbb{E}_\pi[G_t | S_t=s]$$

when we sample all actions according to π

how good it is to be in state s if i am following policy π

Action value function how good it is to take action 'a' from state 's'

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

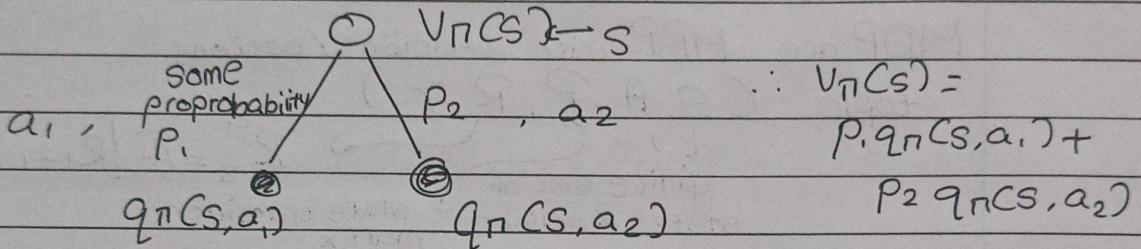
Bellman Expectation Eqn

$$v_n(s) = \mathbb{E}_n [R_{t+1} + \gamma v_n(s_{t+1}) \mid s_t=s]$$

immediate
reward

$$q_n(s, a) = \mathbb{E} [R_{t+1} + \gamma q_n(s_{t+1}, a_{t+1}) \mid s_t=s, a_t=a]$$

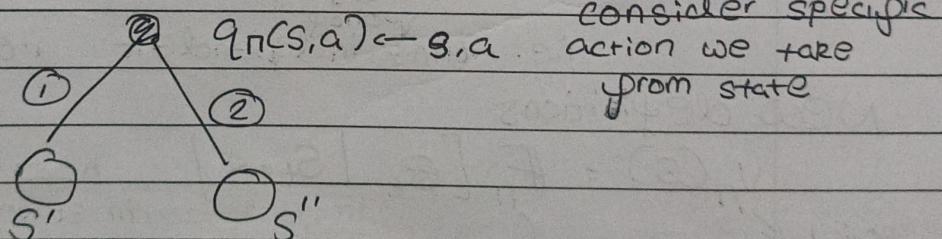
v & q relate to each other



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

Converse relation

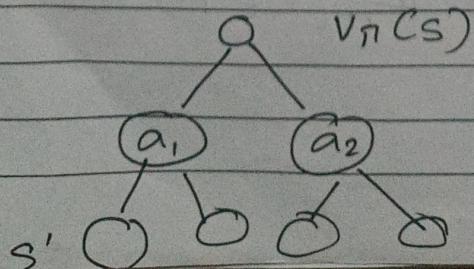
Environment
might blow
you over ① or ②
So $v_n(s)$ gives
how good ① &
② are



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

Combining both

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



MDP can be converted to MRP

& other

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

Basic idea

value function = immediate reward +
value function of
next state

Optimal value function (best behaviour/in MDP)
solution

to pick best policy, see which policy gives max value

$$V^*(s) = \max v_n(s)$$

↳ max possible reward to extract from system

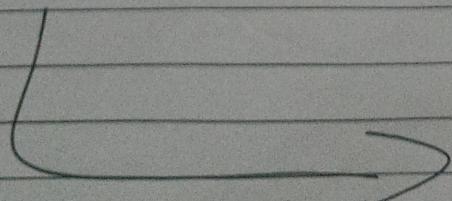
Similarly optimal action value max greward you
get when you fix state & action

$$\underbrace{q^*(s, a)}_{\downarrow} = \max q_n(s, a)$$

Sufficient for RL model

MDP is "solved" when you know q^*

Optimal Policy



Optimal Policy

$$\pi \geq \pi' \text{ if } v_{\pi}(s) \geq v_{\pi'}(s) \text{ for all states}$$

All MDP have ~~at least~~ unique optimal policy (could be more than 1)

$$\pi_* \geq \pi \quad \forall \pi$$

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

optimal value function

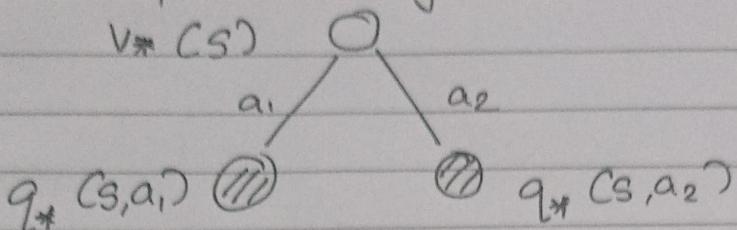
$$q_{\pi_*}(s, a) = q_*(s, a)$$

$$\pi_*(s) = \begin{cases} 1 & \text{if } a = \underset{a \in A}{\operatorname{argmax}} q_*(s, a) \\ 0 & \text{otherwise} \end{cases}$$

Pick action with more q

how to arrive at q_* : Bellman Optimality

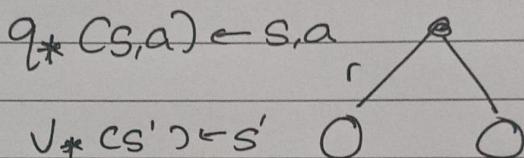
Bellman optimal Policy

(how to figure out q^*)

how good is it to be in state s , is $= \max_{q^*}$ value of actions you can perform

$$V^*(s) = \max_a q^*(s, a) \quad (i)$$

Opp



WE DONT KNOW FOR CERTAIN WHERE ENVIRONMENT MIGHT TAKE US

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

(Combining (i) & (ii) we get

(2 step look ahead)

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

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

This is not linear $\left(\max_{a'} \right)$

→ value iteration

→ policy iteration

→ Q-Learning

→ SARSA

so methods to solve