


MARKOV DECISION PROCESSES

- Planning in uncertain domain.

MDP: It is a discrete time state transition system. It has 4 components

1) S : states.

Ex: For a robot the state could be room, or the (x, y) position.

- play the role of outcomes in the game with environment viewpoint.

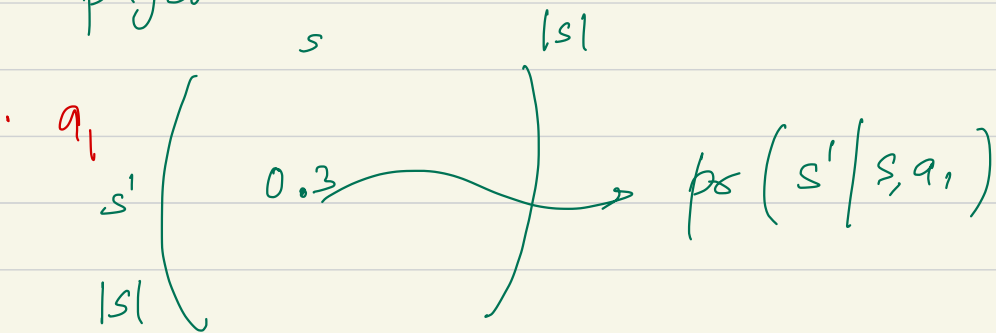
2) A : actions

Usually chosen from a finite set.

3) $\Pr(S_{t+1} | S_t, a_t)$: transition probabilities

. describe the dynamics.

$\phi(s'|s, a)$: the probability of going to state s' from state s if action a is played.



⋮



- By construction next state only depends upon current state and action.

(4) Reward function on states.

$$R: S \rightarrow \mathbb{R}$$

HISTORY: tuples $T_t = (s_0, a_0, s_1, a_1, \dots)$

- Sometimes we specify that s_0 is sampled from a distribution μ on states.

WRITE :

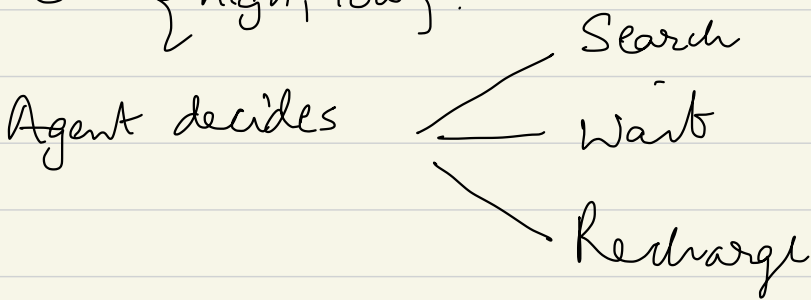
$$M = (S, A, P, R, \gamma, \mu)$$

Example from Barto Sutton.

Recycling robot:

- Collects empty cans.
- Decisions to be made based on current level of "charge"

$S = \{ \text{high, low} \}$.

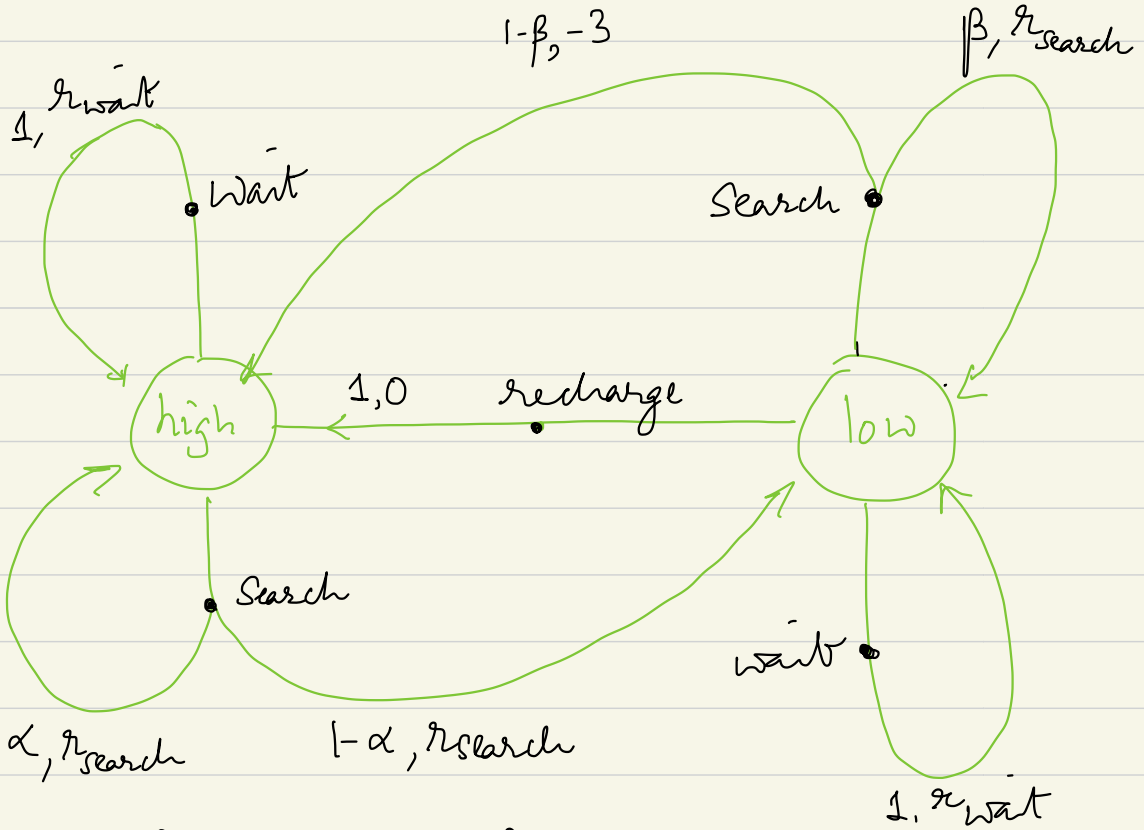


high $\rightarrow \{ \text{search, wait} \}$

low $\rightarrow \{ \text{search, wait, recharge} \}$

Rewards: 0 most times

- positive if a car is found
- negative if battery drains.



- $p(\text{low} | \text{low}, \text{search}) = \beta$
- $r(\text{low}, \text{search}, \text{low}) = r_{search}$

Figure depicts.

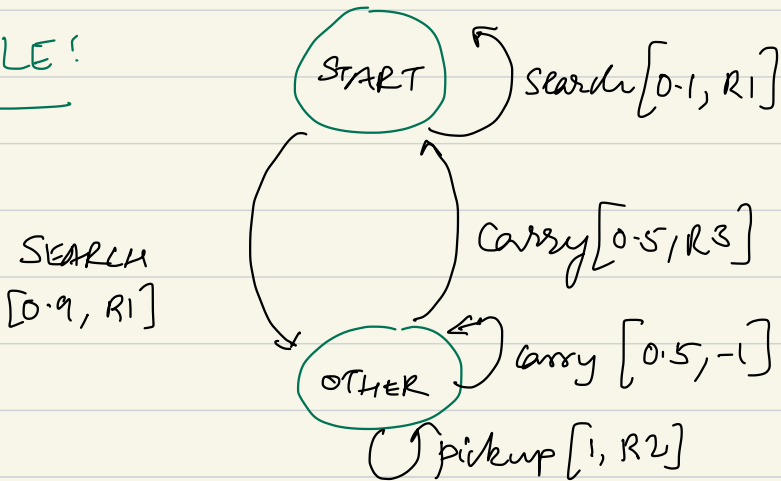
- state nodes



- action nodes •

GOAL: Compute the value function and Q-value.

EXAMPLE:



Def:

a) Policy: $\pi: \mathcal{H} \rightarrow A$ is a map from histories to actions.

b) Deterministic, Stationary policy: is a strategy based on the current state $a_t = \pi(s_t)$

$$\pi: S \rightarrow A$$

c) Stochastic policy:

$$\pi: S \rightarrow \Delta(A)$$

↖ distribution on actions

We write:

$$a_t \sim \pi(\cdot | s_t)$$

the action on the t -th step is drawn from the distribution $\pi(s_t)$

Policy value:

Expected reward when starting at s and following policy π .

Finite horizon:

$$V_{\pi}(s) \stackrel{\text{def}}{=} \mathbb{E}_{a_t \sim \pi(s_t)} \left[\sum_{t=0}^T r(s_t, a_t) \mid s_0 = s \right]$$

The expectation taken over the random selection of actions, and random state s_t reach.

Infinite horizon:

Value function: For a policy π , the value function is the average, discounted sum of future rewards

$$V^\pi(s) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid \pi, s_0 = s \right]$$

Note 1: Assuming $r(s, a)$ is bounded b/w $0 \leq 1$, each term is at most $1 + \gamma + \gamma^2 + \dots$
- bounded by $\frac{1}{1-\gamma}$.

b) So we may normalize it

$$V^\pi(s) = (1-\gamma) \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid \pi, s_0 = s \right]$$

• Note we are working over an infinite horizon.

• for a single time step clearly the action giving max instantaneous reward is best.

• If we fix a horizon k :

• we need to find a policy that gives max expected returns from time 0 to k



Could give us a policy that depends upon time a non-stationary policy.

• Popular to consider infinite horizon.

↓

There is an optimal stationary policy

VALUE FUNCTION for Recycling robot:

$$v^*(h) = \max \left\{ \alpha [r_s + \gamma v^*(h)] + (1-\alpha) [r_s + \gamma v^*(l)], \right. \\ \left. [r_w + \gamma v^*(h)] \right\}$$

$$v^*(l) = \max \left\{ \beta r_s - \gamma(1-\beta) + \gamma [(1-\beta)v^*(h) + \beta v^*(l)] \right. \\ \left. [r_w + \gamma v^*(l)] \right. \\ \left. \gamma v^*(h) \right\}$$

Def: Q-value:

The Q-value function is defined as

$$Q^\pi(s, a) = (1-\gamma) \mathbb{E} \left[\sum \gamma^t r(s_t, a_t) \mid \pi, s_0=s, a_0=a \right]$$

Note that $a_0 = a$ in the above

Goal: Given $s \in S$, find π maximizing

$$V^\pi(s)$$

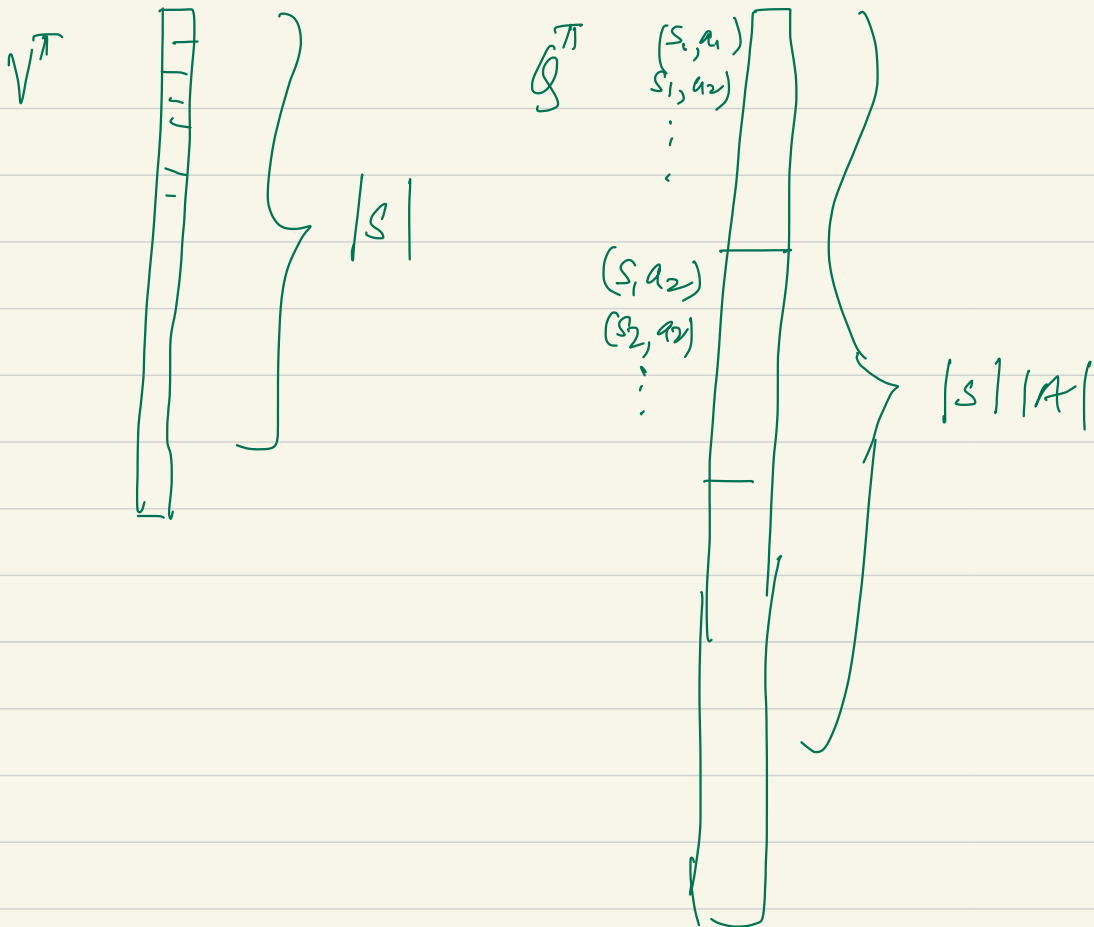
=

By definition: If π is deterministic

$$V^\pi(s) = Q^\pi(s, \pi(s))$$

$$Q^\pi(s, a) = (1-\gamma)r(s, a) + \gamma \mathbb{E}_{s' \sim \mathbb{P}(\cdot | s, a)} [V^\pi(s')]]$$

Called Bellman equation \Leftarrow



• let P be a $|S| |A| \times |S|$ matrix with

$$P_{(s,a), s'} = \mathbb{P}_s(s' | s, a)$$

Define \bar{P}^π , $(|S| \times |A| \times |S| \times |A|)$ matrix

$$\bar{P}^\pi_{(s,a)(s',a')} := \begin{cases} \mathbb{P}_\sigma(s' | s, a) & \text{if } a' = \pi(s') \\ 0 & \text{otherwise} \end{cases}$$

of π 's randomized stationary policy, set

$$\bar{P}^\pi_{(s,a)(s',a')} := \mathbb{P}(s' | s, a) \cdot \pi(a' | s')$$

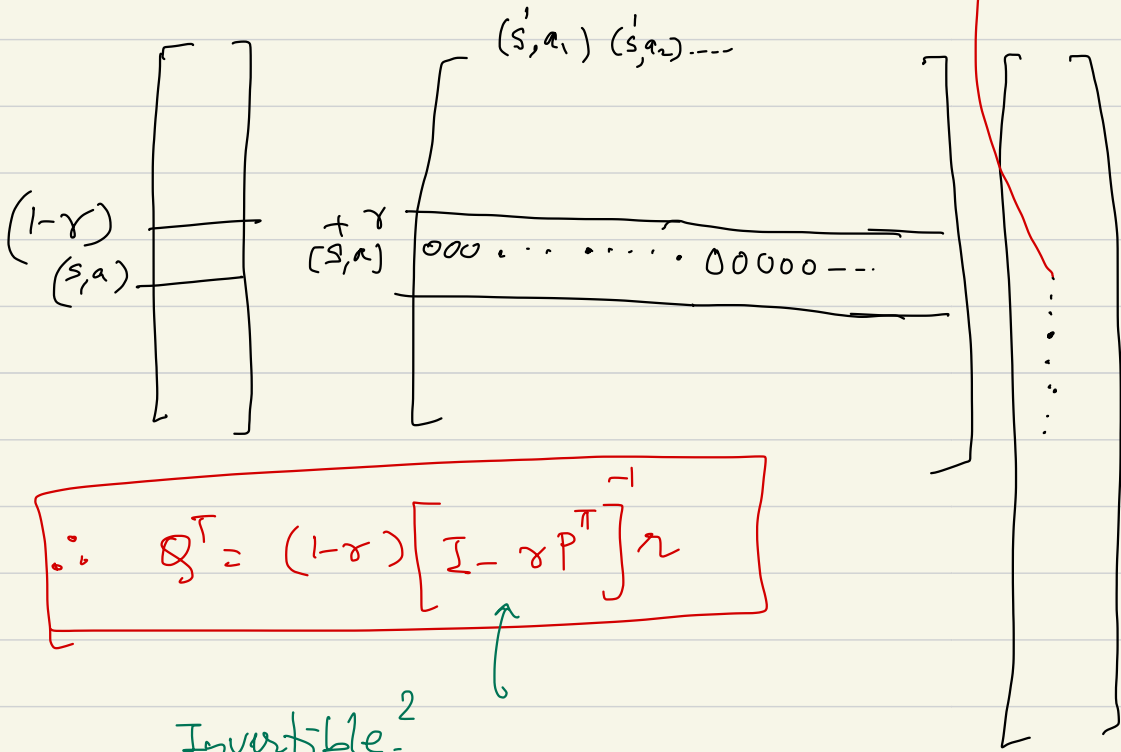
Obs: $Q^\pi = (1-\gamma)r + \gamma \bar{P}^\pi V^\pi$

Pf:
$$\begin{pmatrix} \vdots \\ (s,a) \\ \vdots \end{pmatrix} = (1-\gamma) \begin{pmatrix} \vdots \\ (s,a) \\ \vdots \end{pmatrix} + \gamma \begin{pmatrix} \vdots \\ (s,a) \\ 0 \dots 0 \dots \end{pmatrix} \begin{bmatrix} V^\pi(s_1) \\ \vdots \\ V^\pi(s') \\ \vdots \end{bmatrix}$$

The diagram shows the Bellman optimality equation for the Q-value function. On the left, a vertical vector represents the Q-value for state-action pair (s,a). This is equal to the discounted immediate reward (1-gamma)r plus the discounted expected future value gamma times the sum over next states s' of the transition probability P_sigma(s'|s,a) multiplied by the Q-value for the next state s' and action pi(s'). The transition probability is indicated by an arrow from the state s' in the matrix to the label P_sigma(s'|s,a).

In fact:

$$Q^\pi = (1-\gamma)r + \gamma P^\pi Q$$



Pf: $\|(I - \gamma P^\pi)x\|_\infty = \|x - \gamma P^\pi x\|_\infty$

$x = (x - \gamma P^\pi x) + \gamma P^\pi x$

$\therefore \|x\|_\infty \leq \|x - \gamma P^\pi x\|_\infty + \gamma \|P^\pi x\|_\infty$

$$\begin{aligned}
\therefore \|x - \gamma P^\pi x\|_\infty &\geq \|x\|_\infty - \gamma \|P^\pi x\|_\infty \\
&\geq \|x\|_\infty - \gamma \|x\|_\infty \\
&\geq (1 - \gamma) \|x\|_\infty
\end{aligned}$$

• Bellman optimality equations:

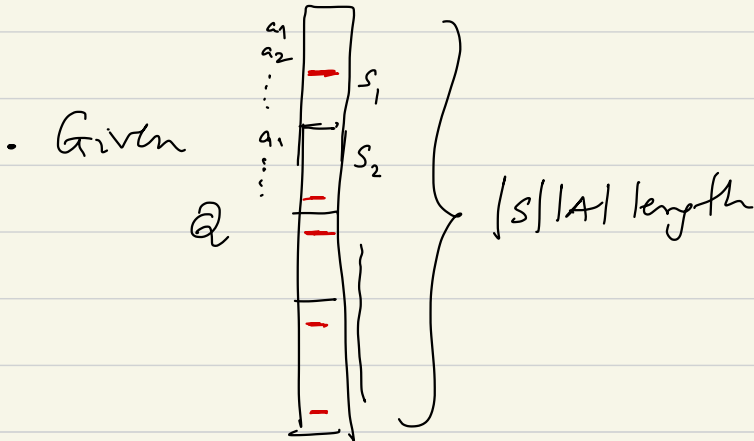
• Thm: For an MDP there is a single stationary and deterministic policy π simultaneously maximizing $V^\pi(s)$ for all $s \in S$ and $Q^\pi(s, a)$, $\forall s \in S, a \in A$;

Denote optimal policy by π^* ;

• V^* corresponding value function
 • Q^* " " Q-value

$$V^*(s) = \max_{a \in A} Q^*(s, a)$$

$$Q^*(s, a) = (1-\gamma)r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} [V^*(s')] \quad *$$



Define:

$$\pi_Q(s) := \operatorname{argmax}_{a \in A} Q(s, a)$$

• So, $\pi^* = \pi_{Q^*}$

• Given Q , define $V_Q(s) := \max_{a \in A} Q(s, a)$

• Bellman operator: $T: \mathbb{R}^{S \times A} \rightarrow \mathbb{R}^{S \times A}$

$$Q \mapsto (1-\gamma)r + \gamma P V_Q$$

* $\implies TQ^* = Q^*$

• Q^* - fixed point for Bellman operator

THM: Let $Q^*(s, a) = \max_{\pi} Q^{\pi}(s, a)$. Then

• \exists a stationary & deterministic policy π s.t. $Q^{\pi} = Q^*$

• $Q \in \mathbb{R}^{S \times A}$ is equal to Q^* iff $TQ = Q$.

= PF:

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

$$= \max_{\pi} \left\{ (1-\gamma) r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} [V^{\pi}(s')] \right\}$$

$$= (1-\gamma) r(s, a) + \gamma \mathbb{E}_{s'} \left[\max_{\pi} V^{\pi}(s') \right] \quad \text{WHY?}$$

$$= (1-\gamma) r(s, a) + \gamma \mathbb{E}_{s'} \left[\max_{\pi} Q^{\pi}(s', \pi(s')) \right]$$

$$= (1-\gamma) r(s, a) + \gamma \mathbb{E}_{s'} \left[\max_{\pi, a'} Q^{\pi}(s', a') \right]$$

" $Q^*(s', a')$

≐ Let $Q = TQ$.

$$\text{Let } \pi = \pi_Q, \therefore Q = (1-\gamma)r + \gamma P^{\pi_Q} Q$$

$$\therefore (1-\gamma)r = (1-\gamma P^{\pi_Q}) Q^{\pi}$$

Let π' be a policy, $Q^{\pi'}$ - the Q -value

$$Q^{\pi'} = (1-\gamma P^{\pi'})^{-1} (1-\gamma)r$$

$$= (1-\gamma P^{\pi'})^{-1} (1-\gamma P^{\pi_Q}) Q^{\pi}$$

$$\therefore Q^{\pi'} - Q^{\pi} = (1-\gamma P^{\pi'})^{-1} \left[(1-\gamma P^{\pi_Q}) - (1-\gamma P^{\pi'}) \right] Q^{\pi}$$

$$= \gamma (1-\gamma P^{\pi'})^{-1} \left[P^{\pi'} - P^{\pi_Q} \right] Q^{\pi}$$

$$\mathbb{E}_{s' \sim P(\cdot | s, a)} \left[\hat{Q}^{\pi}(s', \pi'(s')) - \hat{Q}^{\pi}(s', \pi(s')) \right] \leq 0$$

$\therefore Q^{\pi'} \ll Q^{\pi}$ coordinate wise.

PART-2:

- Q-value iteration algorithm
- Start with random Q .
- Apply $Q \leftarrow TQ$

THM: The above algorithm converges;

PF: (a) show that TQ is contracting

(b) Bound the distance b/w V^{π_Q} & V^* for any Q .

Lemma:

$$\|TQ - TQ'\|_{\infty} \leq \gamma \|Q - Q'\|_{\infty}$$

Pf:

$$\forall s, \left| V_Q(s) - V_{Q'}(s) \right| \leq \max_{a \in A} \left| Q(s, a) - Q'(s, a) \right|$$

$$\text{w.l.o.g. } V_Q(s) > V_{Q'}(s);$$

$$\text{Suppose } V_Q(s) = Q(s, a)$$

$$\text{LHS} = Q(s, a) - \max_{a' \in A} Q'(s, a')$$

$$\leq Q(s, a) - Q'(s, a) \leq \max_a \left| Q(s, a) - Q'(s, a) \right|$$

$$\begin{aligned} &= \|TQ - TQ'\|_{\infty} = \left\| \cancel{(1-\gamma)r} + \gamma P V_Q - \left(\cancel{(1-\gamma)r} + \gamma P V_{Q'} \right) \right\|_{\infty} \\ &= \gamma \|P V_Q - P V_{Q'}\|_{\infty} \end{aligned}$$

$$\leq \gamma \|V_g - V_{g'}\|_\infty$$

$$= \gamma \max_s |V_g(s) - V_{g'}(s)|$$

$$\leq \gamma \max_s \max_a |Q(s, a) - g'(s, a)|$$

← Def'n

$$= \gamma \|Q - g'\|_\infty$$

(2) $\forall g \in \mathbb{R}^{S \times A}$

$$V^{\pi_g} \geq V^* - \frac{2 \|Q - g^*\|_\infty}{1 - \gamma} \mathbf{1}$$

→ Pf: $V^*(s) - V^{\pi_g}(s)$

$$= Q^*(s, \pi^*(s)) - Q^{\pi_g}(s, \pi_g(s))$$

$$= Q^*(s, \pi^*(s)) - Q^*(s, \pi_g(s)) + \underbrace{Q^*(s, \pi_g(s)) - Q(s, \pi_g(s))}_{\text{Term 2}}$$

$$= Q^*(s, \pi^*(s)) - Q^*(s, \pi_g(s)) +$$

$$\gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \left[V^*(s') - V^{\pi_g}(s') \right]$$

$$\leq \underbrace{Q^*(s, \pi^*(s)) - Q(s, \pi^*(s))}_{\text{Term 1}} + \underbrace{Q(s, \pi_g(s)) - Q^*(s, \pi_g(s))}_{\text{Term 2}}$$

$\begin{pmatrix} Q(s, \pi_g(s)) \\ \geq Q(s, \pi^*(s)) \end{pmatrix}$

$$\leq 2 \|Q - Q^*\|_\infty + \gamma \|V^* - V^{\pi_g}\|_\infty$$

$$\Rightarrow \begin{matrix} \boxed{V^*} \\ \boxed{V^{\pi_g}} \end{matrix} - \begin{matrix} \boxed{V^*} \\ \boxed{V^{\pi_g}} \end{matrix} \leq 2 \|Q - Q^*\|_\infty + \gamma \begin{matrix} \boxed{V^*} \\ \boxed{V^{\pi_g}} \end{matrix} - \gamma \begin{matrix} \boxed{V^*} \\ \boxed{V^{\pi_g}} \end{matrix}$$

$$\therefore V^T \geq V^* - \frac{2 \|g - g^*\|_\infty}{1-\gamma} \cdot \mathbb{1}$$

Thm: If $k \geq \frac{1}{1-\gamma} \log\left(\frac{2}{\epsilon(1-\gamma)}\right)$ the Q -value

update algorithm converges to a value vector

$$\tilde{V} \geq V^* - \epsilon \cdot \mathbb{1}.$$

Proof: Start with $Q^0 = 0$, $Q^{k+1} = T Q^k$
 $= T^k Q^0$

Note: $Q^* = T Q^*$

$$\begin{aligned} \|Q^{(k)} - Q^*\|_\infty &= \|T^k Q^0 - T^k Q^*\|_\infty \\ &\leq \gamma^k \|Q^0 - Q^*\|_\infty \leq (1 - (1-\gamma))^k \|Q^0 - Q^*\|_\infty \\ &\leq \exp(- (1-\gamma)k) \end{aligned}$$

NEXT TIME : POLICY UPDATE