

THEORY

1) : - (6)

$$\mathcal{S} = \{s_1, s_2, s_3, s_4, s_5\}$$

terminal state

Given trajectories:

$$T_1: \{s_2, -1, s_3, 2, s_2, 5, s_1, 0, s_1, 3, s_4, 7, s_0\}$$

$$T_2: \{s_4, 3, s_2, -3, s_1, 2, s_1, 5, s_0\}$$

$$T_3: \{s_3, 2, s_3, 1, s_4, 4, s_0\}$$

$$T_4: \{s_1, 1, s_2, 0, s_3, 2, s_4, 2, s_0\}$$

$$T_5: \{s_4, 6, s_1, -1, s_2, -2, s_2, 2, s_3, 6, s_4, 3, s_0\}$$

(any $s_t^{(0,1)}$ will be accepted) (let $t=1$) - as all reach terminal stat - simply add rewards.

(a) First Visit Method:

$$v(s_0) = 0 \text{ (terminal state)}$$

Collecting trajectory rewards :-

$\boxed{\square} \rightarrow$ (sum of reward of a trajectory)

	s_1	s_2	s_3	s_4
T_1	(0, 3, 7)	(-1, 2, 5, 0, 3, 7)	(2, 5, 0, 3, 7)	(7)
T_2	(2, 5)	(-3, 2, 5)	-	(3, -3, 2, 5)
T_3	-	-	(2, 1, 4)	(4)
T_4	(1, 0, 2, 2)	(0, 2, 2)	(2, 2)	(2)
T_5	(-1, -2, 2, 6, 3)	(-2, 2, 6, 3)	(6, 3)	(6, -1, 2, 2, 6, 3)

$$\hat{V} = \frac{1}{30} \begin{pmatrix} 30 \\ 33 \\ 37 \\ 34 \end{pmatrix}$$

$$\hat{V} = \left(\begin{array}{c} \frac{30}{4} \\ 7.5 \\ \hat{V}(1) \end{array} \right) \quad \frac{33}{4} \quad \frac{37}{4} \quad \frac{34}{5}$$

$$8.25 \quad 9.25 \quad 6.8$$

$$\hat{V}(2) \quad \hat{V}(3) \quad \hat{V}(4)$$

(b) Every visit Method:

$$s_1 \quad s_2 \quad s_3 \quad s_4$$

$$T_1 \quad (0, 3, 7)_{10} \quad (-1, 2, 5, 0, 3, 7)_{16} \quad (2, 5, 0, 3, 7)_{17} \quad (7)_{7}$$

$$(3, 7)_{10} \quad (5, 0, 3, 7)_{15}$$

$$T_2 \quad (2, 5)_{7} \quad (-3, 2, 5)_{4} \quad - \quad (3, -3, 2, 5)_{7}$$

$$(5)_{5}$$

$$T_3 \quad - \quad - \quad (2, 1, 4)_{7} \quad (4)_{4}$$

$$(1, 4)_{5}$$

$$T_4 \quad (1, 0, 2, 2)_{5} \quad (0, 2, 2)_{4} \quad (2, 2)_{4} \quad (2)_{2}$$

$$T_5 \quad (-1, -2, 2, 6, 3)_{8} \quad (-2, 2, 6, 3)_{9} \quad (6, 3)_{9} \quad (6, 1, -2, 2, 1, 3)_{14}$$

$$(2, 6, 3)_{11} \quad (3)_{3}$$

$$\begin{array}{cccc} \uparrow : & 40 & \underline{59} & \underline{42} \\ \xrightarrow{* \text{ common}} & \underline{5} & \underline{6} & \underline{7} \\ 8 & 9.83 & 8.4 & 6 \end{array}$$

$$\begin{array}{cccc}
 \text{(* 0m)} & \text{5} & & \\
 8 & \text{9.83} & 8.4 & 6 \\
 \parallel & \parallel & \parallel & \parallel \\
 \hat{V}(1) & \hat{V}(2) & \hat{V}(3) & \hat{V}(4)
 \end{array}$$

2) :- ⑨

\leftarrow (stochastic end time)

$$(a) Q(s, a) = \mathbb{E} \left[\sum_{i=0}^T r(s_i, a_i, s_{i+1}) \mid S_0 = s, A_0 = a \right]$$

$$= \mathbb{E}_{\substack{s' \sim p(\cdot | s, a)}} \left[r(s, a, s') + \sum_{i=1}^T r(s_i, a_i, s_{i+1}) \mid \begin{matrix} S_0 = s \\ A_0 = a \end{matrix} \right]$$

(Controlled ↑

MC)

$$= \mathbb{E} \left[r(s, a, s_1) + \mathbb{E} \left[\sum_{i=1}^T r(s_i, a_i, s_{i+1}) \mid \begin{matrix} s_1 = s' \\ \dots \\ S_0 = s \\ A_0 = a \end{matrix} \right] \right]$$

$$Q(s, a) = \mathbb{E} \left[r(s, a, s') + V(s') \mid S_0 = s, A_0 = a \right]$$

Let V^* and Q^* be the optimal value functions —

So,

$$Q^*(s, a) = \mathbb{E} \left[r(s, a, s') + V^*(s') \mid S_0 = s, A_0 = a \right]$$

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

Therefore, the Bellman optimality eqn:

$$Q^*(s, a) = \mathbb{E}_{\substack{s' \sim p(\cdot | s, a)}} \left[r(s, a, s') + \max_{a' \in A(s')} Q^*(s', a') \mid \begin{matrix} S_0 = s \\ A_0 = a \end{matrix} \right]$$

Thus: the Bellman optimality operators (as discussed in class).

Thus; the Bellman optimality operator's (as discussed in class);

natural choice would be:

$$H : \mathbb{R}^{|S| \times |A|} \rightarrow \mathbb{R}^{|S| \times |A|}$$

$$H Q(s, a) = \mathbb{E} \left[\tau(s, a, s') + \max_{a' \in A(s')} Q(s', a') \mid s_a = s, a = a \right]$$

$$= \sum_{s' \in S} P(s'|s, a) \left[\tau(s, a, s') + \max_{a' \in A(s')} Q(s', a') \right]$$

[BED]

(b) To show: for proper policy $\exists z \in \mathbb{R}^{|S|}$; $z(s) > 0$

$$\text{s.t. } \sum_{s' \in S} P_\pi(s'|s) z(s') \leq \beta z(s) \text{ for some } \beta \in (0, 1)$$

$$\text{let } v_{\pi_*}(s) = \max_a Q_{\pi_*}(s, a)$$

$$= \max_a \mathbb{E} \left[\tau(s, a, s') + v_{\pi_*}(s') \mid s, a \right]$$

$$= \max_{a \in A(s)} \sum_{s' \in S} P(s'|s, a) \left[\tau(s, a, s') + v_{\pi_*}(s') \right]$$

(Is valid for any choice of τ)

safely assume 1

$$= 1 + \max_{a \in A(s)} \sum_{s' \in S} P(s'|s, a) v_{\pi_*}(s')$$

$$\geq \sum_{s' \in S} P(s'|s, a) v_{\pi_*}(s')$$

$a \sim \pi(\cdot | s)$ (for all policies)

by definition

$$> 1 + \sum_{a \sim \pi(s)} \pi(a | s) \sum_{s' \in S} P(s'|s, a) v_{\pi_*}(s')$$

$$= 1 + \underbrace{\sum_{s'} \left(\sum_{a'} p(s'|s, a) \pi(a|s) \right) v_{\pi}(s')}_{P_{\pi}(s'|s)}$$

$$v_{\pi}(s) \geq 1 + \sum_{s'} P_{\pi}(s'|s) v_{\pi}(s')$$

Def; $Z(s) = v_{\pi}(s)$,

$$\Rightarrow \sum_{s' \in S} P_{\pi}(s'|s) Z(s') \leq Z(s) - 1$$

Def $\beta = \max_s \left(\frac{Z(s) - 1}{Z(s)} \right) \in (0, 1)$

$$\Rightarrow \sum_{s' \in S} P_{\pi}(s'|s) Z(s') \leq \beta Z(s)$$

$\because (\exists \beta \text{ for } Z)$

(QED)

(C) To show \otimes is a contraction under a Z norm :-

Let \otimes be a vector of dimension $|S| \times |A|$ such that:

$$\otimes = (\otimes(s, a); \forall s \in S, a \in A(s))^T$$

Consider two \otimes -value functions \otimes , $\bar{\otimes}$

from (a),

$$\therefore H\otimes(s, a) = \sum_{s' \in S} p(s'|s, a) [\otimes(s, a, s') + \max_{a' \in A(s')} \otimes(s', a')]$$

$$\Rightarrow H\otimes(s, a) - H\bar{\otimes}(s, a) = \sum_{s' \in S} p(s'|s, a) \left[\max_{a' \in A(s')} |\otimes(s', a') - \bar{\otimes}(s', a')| \right]$$

(absolute value)

From lemma $\Rightarrow \leq \sum_{s' \in S} p(s'|s, a) \max_{a' \in A(s')} |\otimes(s', a') - \bar{\otimes}(s', a')|$

$$\text{from lemma} \rightarrow \leq \sum_{s' \in S} p(s'|s,a) \max_{a' \in A(s')} |\mathcal{Q}(s',a') - \bar{\mathcal{Q}}(s',a')| \quad (+)$$

Lemma * $\max_x f(x) - \max_a g(x) \leq \max_x |f(x) - g(x)|$

Proof:

$$f(x) - g(x) \leq \max_x (f(x) - g(x))$$

$$\Rightarrow f(x) - \max_x g(x) \leq \max_x (f(x) - g(x))$$

$$\Rightarrow f(x) - \max_x g(x) \leq \max_x (f(x) - g(x))$$

(true for all x)

$$\therefore \max_x f(x) - \max_a g(x) \leq \max_x (f(x) - g(x))$$

$$\leq \max_x |f(x) - g(x)|$$

(definition of absolute value)

Hence

$$H\mathcal{Q}(s,a) - H\bar{\mathcal{Q}}(s,a) \leq \sum_{s' \in S} p(s'|s,a) \max_{a' \in A(s')} |\mathcal{Q}(s',a') - \bar{\mathcal{Q}}(s',a')|$$

Note :

(we can put 1.1 on LHS as if we replace \mathcal{Q} by $\bar{\mathcal{Q}}$,
the RHS of eqn 1 still remains the same)

multipy and divide by $Z(s)$

$$\therefore |H\mathcal{Q}(s,a) - H\bar{\mathcal{Q}}(s,a)| \leq \sum_{s' \in S} p(s'|s,a) \max_{a' \in A(s')} |\mathcal{Q}(s',a') - \bar{\mathcal{Q}}(s',a')| Z(s')$$

$$\leq \sum_{s' \in S} p(s'|s,a) Z(s') \max_{a' \in A(s')} |\mathcal{Q}(s',a') - \bar{\mathcal{Q}}(s',a')|$$

\uparrow
(added)

$$= \| \mathbb{H}\mathbb{Q} - \bar{\mathbb{Q}} \|_2 \sum_{S' \in S} P(S'|S, a) \bar{z}(S')$$

From theorem in (b), $\exists \beta \in [0, 1]$ s.t.

$$| \mathbb{H}\mathbb{Q}(S, a) - \mathbb{H}\bar{\mathbb{Q}}(S, a) | \leq \| \mathbb{H}\mathbb{Q} - \bar{\mathbb{Q}} \|_2 \beta z(S)$$

$$\Rightarrow \max_{S, a} \frac{| \mathbb{H}\mathbb{Q}(S, a) - \mathbb{H}\bar{\mathbb{Q}}(S, a) |}{z(S)} \leq \beta \| \mathbb{H}\mathbb{Q} - \bar{\mathbb{Q}} \|_2 \quad (\text{true for all } a, S)$$

$$\Rightarrow \| \mathbb{H}\mathbb{Q} - \mathbb{H}\bar{\mathbb{Q}} \|_2 \leq \beta \| \mathbb{Q} - \bar{\mathbb{Q}} \|_2.$$

\therefore (H is a contraction under z norm)
 (QED)

(d) Now; since H is a contraction and.

$$\mathbb{H}\mathbb{Q}^* = \mathbb{Q}^* \quad (\text{Bellman optimality eqn})$$

By Banach's fixed pt theorem; there exists a unique fixed pt for H . and \mathbb{Q}^* is the fixed point.

\therefore \mathbb{Q} -learning algorithm using H converges to \mathbb{Q}^*

Thus; we proved convergence of $\overset{H}{\mathbb{Q}}$ -learning algorithm.
 (tabular and model based)

3) : - Q1

(a) Definition:

$$V_{\pi}(s) = \mathbb{E} \left[\sum_{t=0}^{\infty} r^t \mathbb{P}(s_{t+1}, a_{t+1}, s_{t+2}, \dots) \mid s_t = s \right]$$

$$V_{\pi}(s) = \mathbb{E}_{\substack{a \sim \pi(\cdot | s) \\ s' \sim P(\cdot | s, a)}} \left[\sum_{i=0}^{\infty} \gamma^i r(s_{t+i}, a_{t+i}, s_{t+i+1}) \mid s_t = s \right]$$

$$= \mathbb{E} \left[r(s_t, a_t, s_{t+1}) + \gamma r(s_{t+1}, a_{t+1}, s_{t+2}) \right. \\ \left. + \sum_{i=2}^{\infty} \gamma^i r(s_{t+i}, a_{t+i}, s_{t+i+1}) \mid s_t = s \right]$$

$$= \mathbb{E} \left[r(s_t, a_t, s_{t+1}) + \gamma r(s_{t+1}, a_{t+1}, s_{t+2}) \right. \\ \left. + \gamma^2 \mathbb{E}_{\substack{t=0 \\ \pi}} \left[\sum_{i=0}^{\infty} \gamma^i r(s_{t+i+2}, a_{t+i+2}, s_{t+i+3}) \mid s_{t+2} \right] \right. \\ \left. \left. \mid s_t = s \right] \right]$$

(iterated Expectation)

$$= \mathbb{E} \left[r(s_t, a_t, s_{t+1}) + \gamma r(s_{t+1}, a_{t+1}, s_{t+2}) + \gamma^2 V_{\pi}(s_{t+2}) \mid s_t = s \right]$$

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

$s_{t+1} \sim P(\cdot | s_t, a_t)$

(2 step reward function)

$a_{t+1} \sim \pi(\cdot | s_{t+1})$

$s_{t+2} \sim P(\cdot | s_{t+1}, a_{t+1})$

Hence proved

(b)

Objective - linear

$$d(\omega) = \frac{1}{2} \sum_{s \in S} d_{\pi}(s) (V_{\pi}(s) - \Phi^T(s) \omega)^2$$

Gradient descent based update:

$$\frac{\partial L}{\partial \omega} = - \sum_{s \in S} d_{\pi}(s) (V_{\pi}(s) - \Phi^T(s) \omega) \Phi(s) \\ = - \mathbb{E}_{s \sim d_{\pi}} [(V_{\pi}(s) - \Phi^T(s) \omega) \Phi(s)]$$

From part (a),

$$= \mathbb{E}_{s \sim d_{\pi}} \left[\mathbb{E} \left[r(s_t, a_t, s_{t+1}) + \gamma r(s_{t+1}, a_{t+1}, s_{t+2}) \right. \right. \\ \left. \left. + \gamma^2 V_{\pi}(s_{t+2}) - \Phi^T(s_t) \omega \right] \right]$$

$$+ \gamma^2 v_{\pi}(s_{t+2}) - \Phi^T(s_t) w \Big] \\ \Phi(s) \Big)$$

$$= \mathbb{E}_{\substack{s \sim d_{\pi}(\cdot) \\ a_t \sim \pi(\cdot | s_t) \\ s_{t+1} \sim P(\cdot | s_t, a_t) \\ a_{t+1} \sim \pi(\cdot | s_{t+1}) \\ s_{t+2} \sim P(\cdot | s_{t+1}, a_{t+1})}} \left[\left(r(s_t, a_t, s_{t+1}) + \gamma r(s_{t+1}, a_{t+1}, s_{t+2}) + \gamma^2 v_{\pi}(s_{t+2}) \right) - \Phi^T(s_t) w \right] \Phi(s_t)$$

: SGD update: replace the \mathbb{E} expectation and update per sample.

$$w_{t+1} = w_t + \alpha_t' \left(r(s_t, a_t, s_{t+1}) + \gamma r(s_{t+1}, a_{t+1}, s_{t+2}) + \gamma^2 v_{\pi}(s_{t+2}) \right. \\ \left. - \Phi^T(s_t) w_t \right] \Phi(s_t)$$

$v_{\pi}(s_{t+2})$ is unknown: (Sutton's Idea)

$$\text{def } v_{\pi}(s_{t+2}) \approx \Phi^T(s_{t+2}) w_t$$

Algorithm:

$$\Rightarrow w_{t+1} = w_t + \alpha_t' \left(r(s_t, a_t, s_{t+1}) + \gamma r(s_{t+1}, a_{t+1}, s_{t+2}) + \gamma \Phi^T(s_{t+2}) w_t \right. \\ \left. - \Phi^T(s_t) w_t \right] \Phi(s_t)$$

Iterate until convergence - with samples s_{t+1} 's based on $P(\cdot | s_t, a_t)$

and a_t 's based on π .

(c) limiting ODE:

def:

$$f(w_t, s_t, a_t, s_{t+1}, a_{t+1}, s_{t+2})$$

$$= \left(r(s_t, a_t, s_{t+1}) + \gamma r(s_{t+1}, a_{t+1}, s_{t+2}) + \gamma^2 \Phi^T(s_{t+2}) w_t - \Phi^T(s_t) w_t \right) \\ \Phi(s_t)$$

let $f_t = \sigma(w_0, s_0, a_0, s_1, w_1, s_1, a_1, \dots, w_t, s_t)$ be a

Collection of observable for those s.v.s

$$\therefore f(w_t, s_r) = \mathbb{E} [f(w_r, s_r, a_r, s_{r+1}, a_{r+1}, s_{r+2}) \mid f]$$

$$= \mathbb{E} \left[\underbrace{\left(r(s_r, a_r, s_{r+1}) + \gamma r(s_{r+1}, a_{r+1}, s_{r+2}) + \gamma^2 \Phi^T(s_{r+2}) w_r - \Phi^T(s_r) w_r \right)}_{\text{known}} \mid w_r, s_r \right]$$

) expected reward.

$$= r(s_r) \Phi(s_r) - \Phi(s_r) \Phi(s_r) w_r + \gamma \mathbb{E} \left[\underbrace{r(s_{r+1}, a_{r+1}, s_{r+2})}_{\text{III}} + \gamma \Phi^T(s_{r+2}) w_r \mid w_r, s_r \right] \Phi(s_r) \quad (C-1)$$

$$\text{III} = \mathbb{E} \left[\underbrace{r(s_{r+1}, a_{r+1}, s_{r+2}) + \gamma \Phi^T(s_{r+2}) w_r}_{\text{(iterative exp)}} \mid s_r, w_r \right]$$

\downarrow
(s should be a function of s_r, w_r only)

$$= \gamma \Phi(s_r) \sum_{a_r, s_{r+1}} \pi(a_r | s_r) P(s_{r+1} | s_r, a_r) \mathbb{E} \left[r(s_{r+1}, a_{r+1}, s_{r+2}) + \gamma \Phi^T(s_{r+2}) w_r \mid s_{r+1} \right]$$

$$= \gamma \Phi(s_r) \sum_{a_r, s_{r+1}} \pi(a_r | s_r) P(s_{r+1} | s_r, a_r) \left[\sum_{a_{r+1}, s_{r+2}} \pi(a_{r+1} | s_{r+1}) P(s_{r+2} | s_{r+1}, a_{r+1}) \right. \\ \left. \underbrace{(r(s_{r+1}, a_{r+1}, s_{r+2}) + \gamma \Phi^T(s_{r+2}) w_r)}_{\gamma r(s_{r+1})} \right]$$

$$= \gamma \Phi(s_r) \sum_{s_{r+1}} P(s_{r+1} | s_r) \gamma r(s_{r+1}) + \gamma^2 \Phi(s_r) \sum_{s_{r+1}} P(s_{r+1} | s_r) \sum_{s_{r+2}} \underbrace{P(s_{r+2} | s_{r+1})}_{\pi(s_{r+2} | s_{r+1})} \underbrace{\Phi^T(s_{r+2}) w_r}_{\Phi^T(s_{r+2}) w_r}$$

$$= \left[\sum_{s_{r+2}} \sum_{s_{r+1}} P(s_{r+1} | s_r) P(s_{r+2} | s_{r+1}) \right] \underbrace{\Phi^T(s_{r+2}) w_r}_{w_r}$$

$$\underbrace{P^2(s_{t+2}|s_t)}_{\Phi^T(s_{t+2})w_t}$$

$$= \gamma \Phi(s_r) \sum_{s_{t+1}} P_\pi(s_{t+1}|s_r) \pi(s_{t+1}) + \gamma^2 \Phi(s_r) \sum_{s_{t+2}} P_\pi^2(s_{t+2}|s_r) \Phi^T(s_{t+2}) w_t$$

\therefore Putting in C.1:

$F(w_t, s_r)$

$$= \Phi(s_r) \pi(s_r) + \gamma \Phi(s_r) \sum_{s_{t+1}} P_\pi(s_{t+1}|s_r) \pi(s_{t+1}) + \gamma^2 \Phi(s_r) \sum_{s_{t+2}} P_\pi^2(s_{t+2}|s_r) \Phi^T(s_{t+2}) w_t - \Phi(s_r) \Phi^T(s_r) w_t$$

— (C.2)

Now,

$$\underline{w_{t+1}} = w_t + \alpha_t f(w_t, s_r, a_t, s_{t+1}, a_{t+1}, s_{t+2})$$

$$= w_t + \alpha_t [F(w_t, s_r) + M_{t+1}]$$

$$\underline{\text{where: } M_{t+1} = f(w_t, s_r, a_t, s_{t+1}, a_{t+1}, s_{t+2}) - F(w_t, s_r)}$$

Show M_{t+1} is a martingale:

$$\mathbb{E}[M_{t+1} | \mathcal{F}_t] = \underbrace{\mathbb{E}[f | \mathcal{F}_t]}_{\text{constant given } w_t, s_r} - \mathbb{E}[F(w_t, s_r) | w_t, s_r]$$

$$= F(w_t, s_r) - F(w_t, s_r)$$

(from C.2)

$$= 0 \quad (\text{Hence } M_{t+1} \text{ is a valid Martingale})$$

(d) As seen in class (law of stochastic approximation algorithms); the solution of the SDE is the following ODE:

(d) As seen in class (law of stochastic approximation algorithms); the solution of the SDE is the following ODE:

$$\dot{w}(t) = \sum_{s_t} d_{\text{TF}}(s_t) \text{TF}(w_t, s_t)$$

(Pretty C-2)

$$= \sum_{s_t} d_{\pi}(s_t) \Phi(s_t) \mathcal{J}_2(s_t) + \gamma \sum_{s_{t+1}} d_{\pi}(s_{t+1}) \Phi(s_{t+1}) \sum_{\pi} P_{\pi}(s_{t+1}, s_t) \mathcal{J}_2(s_{t+1})$$

$$+ \gamma^2 \sum_{t=1}^T \underbrace{\Phi(s_t)}_{\sum_{s_{t+1}} p^2(s_{t+1}|s_t)} \underbrace{\Phi^T(s_{t+2})}_{w_t}$$

$$\Phi^T D_{II} \Phi w_f - \sum_s d_{II}(s_f) \Phi(s_f) \Phi^T(s_f) w_f$$

{ $\Phi^T D_{II} \Phi w_f$ }

(quadratic)

$$\text{III} \quad \delta^2 \sum_t d_{\Pi}(s_t) \Phi(s_t) \sum_t P^2(s_{t+2}|s_t) \Phi^T(s_{t+2}) w_t$$

$$\gamma^2 \Phi^T D_{\pi} P_{\pi}^2 \Phi \underset{S \times S}{\overset{S \times d}{\rightarrow}} w_T \underset{d \times 1}{\rightarrow}$$

$$\therefore \hat{w}(t) = \Phi^T D_{\pi} \gamma + \gamma \Phi^T D_{\pi} P_{\pi} \gamma + \gamma^2 \Phi^T D_{\pi} P_{\pi}^2 \Phi w_t - \Phi^T D_{\pi} \Phi w_t.$$

$$= \Phi^T D_{\Pi} (I + \gamma P_{\Pi}) \bar{z} - \Phi^T D_{\Pi} (I - \gamma \frac{P_{\Pi}^2}{1 - \gamma^2}) \Phi w_+$$



$$= b - A w_+.$$

$$\tilde{b} = \Phi^T D_{\Pi} (I + \gamma P_{\Pi}) x$$

$$\tilde{A} = \bar{\Phi}^T D_{\pi} (I - \gamma P_{\pi}^2) \bar{\Phi}$$

(e) To show that \tilde{A} is invertible:

$$\tilde{A} = \bar{\Phi}^T D_{\pi} \bar{\Phi} - \gamma \underbrace{\bar{\Phi}^T D_{\pi} P_{\pi}^2 \bar{\Phi}}_M$$

Compute:

$$|x^T M x| = |x^T D_{\pi} P_{\pi}^2 x| = |x^T D_{\pi}^{1/2} D_{\pi}^{1/2} P_{\pi}^2 x|$$

$$\text{(Cauchy-Schwarz)} \leq \|D_{\pi}^{1/2} x\|_2 \|D_{\pi}^{1/2} P_{\pi}^2 x\|_2$$

$$= \|x\|_{D_{\pi}} \left[\|P_{\pi}^2 x\|_{D_{\pi}} \right]^2$$

$$\|P_{\pi}^2 x\|_{D_{\pi}}^2 = \sum_s d_{\pi}(s) \left[\sum_{s'} \underbrace{P(s'|s)}_{\text{(prob distribution - 2 step)}} x(s') \right]^2$$

(Jensen's inequality)

$$\text{Convex: } f(E(x)) \leq E[f(x)] \leq \sum_s d_{\pi}(s) \sum_{s'} P(s'|s) x(s')$$

$$= \sum_{s'} x(s')^2 \sum_s d_{\pi}(s) P(s'|s)$$

$$= \sum_{s'} x(s')^2 d_{\pi}(s')$$

$$= \underbrace{\|x\|_{D_{\pi}}^2}_{\text{Def}} \Rightarrow \|P_{\pi}^2 x\|_{D_{\pi}}^2 \leq \|x\|_{D_{\pi}}^2$$

$$|x^T M x| \leq \|x\|_{D_{\pi}} \|x\|_{D_{\pi}} = \|x\|_{D_{\pi}}^2$$

$$|x^T M x| \leq \|x\|_D \underbrace{\|x\|_{D_\pi}}_{\pi} = \|x\|_D$$

(e.1)

Now,

$$|x^T A x| = \underbrace{x^T \Phi^T D_\pi}_{z} (I - \gamma^2 P_\pi^2) \underbrace{\Phi x}_{z}$$

(let Φ be full column rank: $\Phi x \neq 0$ if $w \neq 0$)

$$= z^T D_\pi (I - \gamma^2 P_\pi^2) z.$$

$$= \|z\|_D^2 - \gamma^2 \|z\|_D^2 \underbrace{\frac{P_\pi^2}{m}}$$

(from e.1)

$$\geq \|z\|_D^2 - \gamma^2 \|z\|_D^2.$$

$$= (1 - \gamma^2) \|z\|_D^2 > 0 \quad (\text{knee positive defn})$$

$\therefore \tilde{A}^{-1}$ exists

$$w_*' = [\tilde{A}^T D_\pi (I - \gamma^2 P_\pi^2) \tilde{A}]^{-1} \underbrace{\tilde{A}^T D_\pi (I + \gamma P_\pi)}_b z$$

\tilde{A}^{-1}

[CASE]

$$\tilde{A} w_* = \tilde{A}^T (\tilde{A}^T \tilde{A})^{-1} \tilde{A}^T D_\pi V_\pi \leftarrow \text{(optimal)}$$

(f) New Bound: Use the following lemma from class.

Lemma 1: $T_\pi V = \pi_2 + \gamma P_\pi V$ is a contraction

Lemma 2: Let $\|v - \pi v\|_D \leq \min_w \|v - \Phi w\|_D$.
 Then $\|\pi v\|_D \leq \|v\|_D$

Lemma 3: πT_π is a contraction w.r.t $\|\cdot\|_D$.

~~We show~~
Lemma 4 $T_\pi^2 v_\pi = T_\pi (\pi v)$ ($\det T_\pi^2$ is a contraction)

$$\begin{aligned} & \|T^\pi(T_\pi v_1 - T_\pi v_2)\|_D \\ & \quad \left(\begin{array}{l} \text{contraction} \\ \text{on the} \end{array} \right) \lesssim \|T_\pi(\gamma + \rho_\pi v_1 - \gamma - \rho_\pi v_2)\| \\ & \quad = \gamma \|P_\pi(T_\pi v_1 - T_\pi v_2)\|_D \\ & \quad \leq \gamma \|T_\pi v_1 - T_\pi v_2\|_D \\ & \quad \text{(as seen in class: } \|P_\pi z\|_D \leq \|z\|_D) \\ & \quad (\text{Lemma 1}) \leq \gamma^2 \|v_1 - v_2\|_D \end{aligned}$$

$(T_\pi^2$ is a γ^2 -contraction)

$$\begin{aligned} \text{Now: } & \|\pi T_\pi^2 v_1 - \pi T_\pi^2 v_2\|_D \\ & = \|\pi(T_\pi^2 v_1 - T_\pi^2 v_2)\|_D \stackrel{(\text{Lemma 2})}{\leq} \|\tilde{T}_\pi v_1 - \tilde{T}_\pi v_2\|_D \\ & \leq \gamma^2 \|v_1 - v_2\|_D. \end{aligned}$$

πT_π^2 is also a γ^2 -contraction.

Now what is $\pi \tau_{\pi}^2$?

$$\pi \tau_{\pi}^2 v = \pi \tau_{\pi} (\tau_{\pi} v)$$

$$= \pi \tau_{\pi} (\gamma + \gamma P_{\pi} v)$$

$$= \pi (\gamma + \gamma P_{\pi} (\gamma + \gamma P_{\pi} v))$$

$$= \pi [(\mathbf{I} + \gamma P_{\pi}) \gamma + \gamma^2 P_{\pi}^2 v]$$

Showing that $\Phi w'_*$ is the fixed pt for $\pi \tau_{\pi}^2$.

$$\pi \tau_{\pi}^2 [\Phi w'_*]$$

$$= \pi [(\mathbf{I} + \gamma P_{\pi}) \gamma + \gamma^2 P_{\pi}^2 \Phi w'_*]$$

$$= \pi (\mathbf{I} + \gamma P_{\pi}) \gamma + \gamma^2 \pi P_{\pi}^2 \Phi w'_* \quad (\text{f.1})$$

[We have seen in class: $\pi = \Phi (\Phi^T D \Phi)^{-1} \Phi^T D$]

$$\pi (\mathbf{I} + \gamma P_{\pi}) \gamma = \Phi (\Phi^T D \Phi)^{-1} \Phi^T D (\mathbf{I} + \gamma P_{\pi}) \gamma$$

(f.2)

Now

form the soln: $A w'_* = b$

we know: (parts c-d)

$$\Rightarrow \Phi^T D (\mathbf{I} - \gamma^2 P_{\pi}^2) \Phi w'_* = \Phi^T D (\mathbf{I} + \gamma P_{\pi}) \gamma$$

$$\Rightarrow \Phi^T D(I - \gamma^2 P_{\pi}^2) \Phi w_*' = \Phi^T D(I + \gamma P) v_2$$

: (f-2 becomes)

$$\begin{aligned} \Pi(I + \gamma P_{\pi}) v_2 &= \cancel{\Phi (\Phi^T D_{\pi} \Phi)^{-1} (\Phi D_{\pi} \Phi)} w_*' \\ &\quad - \cancel{\gamma \Phi (\Phi^T D_{\pi} \Phi)^{-1} \Phi^T D P_{\pi} \Phi} w_*' \\ &= \cancel{\Phi w_*'} - \gamma^2 \Pi P_{\pi}^2 \cancel{\Phi w_*'} \underbrace{\Pi}_{\cancel{\Phi}} \end{aligned}$$

Putting in f-1) $\Pi T_{\pi}^2 \Phi w_*' =$

$$\begin{aligned} &\cancel{\Pi} ((I + \gamma P_{\pi}) v_2 + \gamma^2 P_{\pi}^2 \Phi w_*') \\ &= \cancel{\Phi w_*'} - \cancel{\gamma^2 \Pi P_{\pi} \Phi w_*'} + \cancel{\gamma^2 \Pi P_{\pi}^2 \Phi w_*'} \\ &= \Phi w_*' \end{aligned}$$

$\Phi w_*'$ is the fixed pt of ΠT_{π}^2
(operator of interest)

Now - Bound

$$= \| \Phi w_*' - v_{\pi} \|$$

(find it)

$$= \| \Pi T_{\pi}^2 \Phi w_*' - \Phi w_*' + \Phi w_*' - v_{\pi} \| \quad \text{(optimal)}$$

Triangle inequality

<

$$\| \Pi T_{\pi}^2 \Phi w_*' \| + \| \Phi w_*' - v_{\pi} \|$$

\leq

$$\|\underbrace{\pi \tau_{\pi}^2 \Phi w_*' - \Phi w_{\#}}_D\| + \|\Phi w_{\#} - v_{\pi}\|_D.$$

πv_{π} (optimal sol^r)

$$\|\tau_{\pi} v_{\pi}\|$$

$$\tau_{\pi}^2 v_{\pi} = \tau_{\pi} (\tau_{\pi} v) = v_{\pi}.$$

$$= \|\underbrace{\pi \tau_{\pi}^2 \Phi w_*'}_D - \underbrace{\tau_{\pi}^2 v_{\pi}}_D\| + \|\Phi w_{\#} - v_{\pi}\|_D.$$

(It is a γ^2 contraction - lemma 4)

$$\leq \gamma^2 \|\Phi w_*' - v_{\pi}\|_D + \|\Phi w_{\#} - v_{\pi}\|_D.$$

$$\Rightarrow (1 - \gamma^2) \|\Phi w_*' - v_{\pi}\|_D \leq \|\Phi w_{\#} - v_{\pi}\|_D.$$

$$\|\Phi w_*' - v_{\pi}\|_D \leq \frac{1}{1 - \gamma^2} \|\Phi w_{\#} - v_{\pi}\|_D.$$

$$\text{old bound } \|\Phi w_*' - v_{\pi}\|_D \leq \frac{1}{1 - \gamma} \|\Phi w_{\#} - v_{\pi}\|_D$$

(1 step reward)

The new algorithm is better for a tighter upper bound.

$$\gamma^2 < \gamma \quad * \quad \gamma \in (0, 1)$$

$$1-\gamma^2 > 1-\gamma$$

$$\Rightarrow \frac{1}{1-\gamma^2} < \frac{1}{1-\gamma} \quad (\text{Tighter upper bound
- better algorithm})$$

Coding = (25 pts) — see Assignment2-3.ipynb code.

Q.1 : Model-free TD(γ) algorithm -

Update :

$$\forall s \in S, \quad V_{k+1}(s) = V_k(s) + \left[\gamma \sum_{m=K}^L \alpha_{km} d_m \right] \quad \begin{matrix} \leftarrow \text{(End of episode)} \\ L \\ m-K \end{matrix} \quad \alpha_K = \frac{1}{\# \text{visits}}$$

$$\text{where : } \bar{d}_m = \pi(s_m, a_m, s_{m+1}) + V_K(s_{m+1}) - V_K(s_m)$$

where ($s_{km} = s$)

(γ is the discount factor), $\{\lambda\}_{m=K}^{m=L}$ is the weight
for future errors
 d_m .

Algorithm :

Overall idea :

$$1) \underline{\text{Store}}: \left\{ \begin{array}{l} [s_0, s_1, s_2, \dots, s_L] = S \\ [d_0, d_1, d_2, \dots, d_L] = D \end{array} \right.$$

$$\left([d_0 \ d_1 \ \dots \ \dots \ \dots \ d_L] \right) = D$$

(And at the end of the trajectory)

 (temporal differences)

Do dp:

for $\vartheta \in \{L-1, L-2, \dots, 0\}$

$$\text{Cumu}(D)[\vartheta] = \lambda \text{Cumu}[\vartheta+1] + D[\vartheta].$$



$$V(S[i]) := \text{Cumu}[i] \quad \forall i \in \{0, 1, 2, \dots\}$$

Q.2:

Bellman Equation:

$$Q_{\pi}(s, a) = \sum_{s', a'} P(s'|s, a) \left[r(s, a, s') + \gamma \max_{a' \in A(s')} Q_{\pi}(s', a') \right]$$

Q.3:

$$\omega_{t+1} = \omega_t + \alpha_t \left[r(s_t, a_t, s_{t+1}) + \gamma \Phi^T(s_{t+1}) \omega_t - \Phi^T(s_t) \omega_t \right] / \Phi(s_t)$$