

HW 1

Problem 1

(a)

(1)

$$\max_{\pi} V^{\pi}(s) = \max_{\pi} \sum_a \pi(a|s) \cdot Q^{\pi}(s, a) \leq \max_{\pi} Q^{\pi}(s, \pi(s)) = Q^{*}(s, a)$$

" ↓

$V^{*}(s)$ $\because 0 \leq \pi(a|s) \leq 1$

Suppose $V^{*}(s) < Q^{*}(s, a)$,

$$\rightarrow \max_{\pi} \sum_a \pi(a|s) \cdot Q^{\pi}(s, a) < \max_{\pi} Q^{\pi}(s, a)$$

$$\text{take } a^{*} = \operatorname{argmax}_{a \in A} Q^{\pi}(s, a) \wedge \pi(a^{*}|s) = 1$$

$$\rightarrow \max_{\pi} Q^{\pi}(s, a^{*}) < \max_{\pi} Q^{\pi}(s, a^{*}) \rightarrow 0 < 0 \quad (*)$$

$$\therefore V^{*}(s) = Q^{*}(s, a) \quad \#$$

(2)

$$\max_{\pi} (Q(s, a)) = \max_{\pi} \left\{ R(s, a) + \gamma \cdot \sum_{s'} P_{ss'}^a \cdot \max_{\pi} V^{*}(s') \right\}$$

$$\rightarrow \max_{\pi} (Q(s, a)) = R(s, a) + \gamma \cdot \sum_{s'} P_{ss'}^a \cdot \max_{\pi} V^{*}(s')$$

$$\rightarrow Q^{*}(s, a) = R(s, a) + \gamma \cdot \sum_{s'} P_{ss'}^a \cdot V^{*}(s')$$

(b)

$$\begin{aligned} & \|T^*(Q) - T^*(Q')\|_\infty \\ &= \max_{(s,a)} |(\cancel{R_{s,a}} + \gamma \cdot \sum_{s'} P_{ss'}^a \cdot \max_{a'} Q(s', a')) - (\cancel{R_{s,a}} + \gamma \cdot \sum_{s'} P_{ss'}^a \cdot \max_{a'} Q'(s', a'))| \\ &= \max_{(s,a)} \gamma \cdot \left| \sum_{s'} P_{ss'}^a \cdot (\max_{a'} Q(s', a') - \max_{a'} Q'(s', a')) \right| \\ &\leq \max_{(s,a)} \gamma \cdot \left| \sum_{s'} P_{ss'}^a \cdot \max_{a'} (Q(s', a') - Q'(s', a')) \right| \quad (|\max A - \max B| \leq \max(A-B)) \\ &\leq \gamma \cdot \max_{(s,a)} \max_{a'} |Q(s', a') - Q'(s', a')| \quad (\because 0 \leq \text{entry of } P_{ss'}^a(s,a) < 1) \\ &\leq \gamma \cdot \|Q - Q'\|_\infty \\ &\therefore T^* \text{ is a } \gamma\text{-contraction operator} \quad \# \end{aligned}$$

Problem 1

(a)

$$\begin{aligned} & \|T_\Omega^\pi V(s) - T_\Omega^\pi V'(s)\|_\infty \\ &= \max_s |(\cancel{R_s^\pi} + \Omega(\cancel{\pi(\cdot|s)}) + \gamma \cdot P_{ss'}^\pi \cdot V) - (\cancel{R_s^\pi} + \Omega(\cancel{\pi(\cdot|s)}) + \gamma \cdot P_{ss'}^\pi \cdot V')| \\ &= \max_s |\gamma \cdot P_{ss'}^\pi \cdot (V - V')| \\ &\leq \gamma \cdot \max_s |V - V'| = \gamma \cdot \|V - V'\|_\infty \\ &\therefore T_\Omega^\pi \text{ is } \gamma\text{-contraction operator in } L_\infty \text{ norm} \quad \# \end{aligned}$$

(b)

$$\text{define } [T_{\Omega}^{\pi} V](s) = R_s^{\pi} + \Omega(\pi(\cdot|s)) + \gamma \cdot P_{ss}^{\pi} \cdot V$$

<algorithm>

step 1. initialize $K = 0$ and $V_0(s) = 0 \quad \forall s \in S$

step 2. $V_{K+1} \leftarrow T_{\Omega}^{\pi}(V_K)$

$$\Leftrightarrow V_{K+1}(s) = \max_{a \in A} \{ R_s^{\pi} + \underbrace{\Omega(\pi(\cdot|s))}_{\text{also consider the randomness of } \pi} + \gamma \cdot P_{ss}^{\pi} \cdot V_K \}$$

step 3.

$$Q_{K+1}(s, a) = R_{s,a} + \gamma \cdot \sum_{s' \in S} P_{ss'}^a \cdot V_{K+1}(s')$$

<explain>

in step 2, $\because T_{\Omega}^{\pi}$ is γ -contraction operator, it will converge to final V

\therefore we compute T_{Ω}^{π} iteratively

in step 3, we get Q_{K+1} by equation (7),

$$\text{if } V_{K+1}(s) = V_{\Omega}^{*}(s) \Rightarrow Q_{K+1}(s, a) = Q_{\Omega}^{*}(s, a)$$

#

Problem 3

$$RHS = \frac{1}{1-\gamma} \cdot \mathbb{E}_{s \sim d_{\pi}^{\pi_0}} \left[\mathbb{E}_{a \sim \pi_{\theta}(\cdot|s)} [f(s, a)] \right]$$

$$= \frac{1}{1-\gamma} \cdot \sum_{s_0 \in S} \sum_{s_t \in S} \sum_{a_t \in A} f(s, a) \cdot \mu(s_0) \cdot \pi_{\theta}(a_t|s) \cdot (1-\gamma) \cdot \sum_{\tau=0}^{\infty} \gamma^{\tau} \cdot \mathbb{P}(s_t = s | s_0, \pi)$$

$$= \sum_s \sum_{\tau=0}^{\infty} \gamma^{\tau} \cdot f(s, a) \cdot p_{\mu}^{\pi_0, \tau} \quad (\text{by lemma 1})$$

$$= \mathbb{E}_{\tau \sim p_{\mu}^{\pi_0}} \left[\sum_{\tau=0}^{\infty} \gamma^{\tau} \cdot f(s, a) \right] \quad \#$$

Lemma 1. let trajectory $\tau = (s_0, a_0, s_1, a_1, \dots, s_t, a_t)$

$$p_{\mu}^{\pi}(\tau) = \mu(s_0) \cdot \pi_{\theta}(a_0|s_0) \cdots \pi(a_t|s_t) \cdot \mathbb{P}(s_t | s_{t-1}, a_{t-1}, \pi) \cdot \pi(a_t | s_t)$$

by Markov property $\rightarrow \mu(s_0) \cdot \pi(a_t | s_t) \cdot \prod_{i=1}^t \pi(a_{i-1} | \bigcap_{j=0}^{i-1} s_j, \bigcap_{j=0}^{i-1} a_j) \cdot \mathbb{P}(s_t | \bigcap_{j=0}^{t-1} s_j, \bigcap_{j=0}^{t-1} a_j, \pi)$

by def of conditional probability $\rightarrow \mu(s_0) \cdot \frac{\pi(a_0 | s_0)}{\pi(s_0)} \cdot \frac{\mathbb{P}(a_0, \bigcap_{i=0}^{t-1} s_i, \pi)}{\mathbb{P}(a_0, s_0, \pi)} \cdots \frac{\pi(\bigcap_{i=0}^{t-1} s_i, \bigcap_{i=0}^{t-1} a_i)}{\pi(\bigcap_{i=0}^{t-1} s_i, \bigcap_{i=0}^{t-1} a_i)} \cdot \frac{\mathbb{P}(\bigcap_{i=0}^t s_i, \bigcap_{i=0}^{t-1} a_i, \pi)}{\mathbb{P}(\bigcap_{i=0}^{t-1} s_i, \bigcap_{i=0}^{t-1} a_i, \pi)}$

$$= \mu(s_0) \cdot \pi(a_t | s_t) \cdot \mathbb{P}(\bigcap_{i=0}^t s_i, \bigcap_{i=0}^{t-1} a_i | s_0, \pi)$$

Lemma 2.

$$\sum_{\tau} p_{\mu}^{\pi}(\tau) = \sum_{\tau} \mu(s_0) \cdot \pi(a_t | s_t) \cdot \mathbb{P}(\bigcap_{i=0}^t s_i, \bigcap_{i=0}^{t-1} a_i | s_0, \pi) \quad (\text{by lemma 1})$$

extend τ trajectory $\rightarrow \sum_{s_0} \sum_{s_1} \cdots \sum_{s_t} \sum_{a_0} \cdots \sum_{a_t} \mu(s_0) \cdot \pi(a_t | s_t) \cdot \mathbb{P}(\bigcap_{i=0}^t s_i, \bigcap_{i=0}^{t-1} a_i | s_0, \pi)$

by total probability theorem $\rightarrow \sum_{s_0} \sum_{s_t} \sum_{a_t} \mu(s_0) \cdot \pi(a_t | s_t) \cdot \left[\sum_{s_1} \sum_{s_2} \cdots \sum_{s_{t-1}} \sum_{a_0} \cdots \sum_{a_{t-1}} \mathbb{P}(s_t, \bigcap_{i=0}^{t-1} s_i, \bigcap_{i=0}^{t-1} a_i | s_0, \pi) \right]$

$$= \sum_{s_0} \sum_{s_t} \sum_{a_t} \mu(s_0) \cdot \pi(a_t | s_t) \cdot \mathbb{P}(s_t | s_0, \pi)$$

Problem 4

[illegible]

← This is the result of the PI algorithm
More detail please refer to the code

Problem 5

1. Original dataset : 'maze1d-umaze-v1'

Format:

```
[ 1.0856489  1.9745734  0.00981035  0.02174424]
[ 1.0843927  1.97413   -0.12562364 -0.04433781]
[ 1.0807577  1.9752754 -0.3634883  0.11453988]
...
[ 1.1328583  2.8062387 -4.484303  0.09555068]
[ 1.0883482  2.8068895 -4.4510083  0.06509537]
[ 1.0463258  2.8074222 -4.202244  0.05324839]]
load_datafile: 100%
```

Describe :

it is a $N \times \text{dim}$ -observation array, where each observation is a float from -5 to 2

2. Mujoco dataset : 'halfcheetah-random-v1'

Format :

```
[[-2.8627589e-02 -6.3696302e-02  8.6608730e-02 ... -8.6148446e-03
 -4.0926412e-03  1.5979043e-01]
 [-4.3686085e-02 -2.2352228e-02  1.1845701e-02 ...  1.7004584e+00
 -5.4827838e+00  2.5702784e+00]
 [-8.1463926e-02 -3.5814878e-02 -8.6687684e-02 ... -7.7921929e+00
  1.4037144e+01 -8.9725056e+00]
 ...
 [-5.3456318e-01  3.2453218e+00 -9.1072828e-02 ...  1.3767828e+00
 -9.2018442e+00  1.5579005e+00]
 [-4.9477938e-01  3.1914997e+00 -3.9144486e-01 ... -2.0965147e+00
  2.5790567e+00  1.0007437e+01]
 [-5.1447099e-01  3.2184551e+00 -2.8413078e-01 ... -4.2252550e+00
  9.4621916e+00 -3.1550488e+00]]
load_dataframe: 100%
```

Describe:

it is a $N \times \text{dim}$ -observation array, where each observation is a float ranging from -1 to 10