

ECE 49595 - RL : HW 3



Problem 1:

$$1) P^n(s' | s) = \sum_{a \in A} \pi(a | s) \cdot P(s' | s, a)$$

$$\pi(1) = g, \pi(2) = g, \pi(3) = h$$

$$\left. \begin{aligned} \therefore P^n(1 | 1) &= (1 \times 0.1) + (0 \times 0.8) = 0.1 \\ P^n(2 | 1) &= (1 \times 0.8) + (0 \times 0.1) = 0.8 \\ P^n(3 | 1) &= (1 \times 0.1) + (0 \times 0.1) = 0 \end{aligned} \right\} \text{row 1}$$

$$\left. \begin{aligned} P^n(1 | 2) &= (1 \times 0.1) + (0 \times 0.1) = 0.1 \\ P^n(2 | 2) &= (1 \times 0.1) + (0 \times 0.8) = 0.1 \\ P^n(3 | 2) &= (1 \times 0.8) + (0 \times 0.1) = 0.8 \end{aligned} \right\} \text{row 2}$$

$$\left. \begin{aligned} P^n(1 | 3) &= (0 \times 0.8) + (1 \times 0.1) = 0.1 \\ P^n(2 | 3) &= (0 \times 0.1) + (1 \times 0.1) = 0.1 \\ P^n(3 | 3) &= (0 \times 0.1) + (1 \times 0.8) = 0.8 \end{aligned} \right\} \text{row 3}$$

$$\therefore P^n = \begin{bmatrix} 0.1 & 0.8 & 0.1 \\ 0.1 & 0.1 & 0.8 \\ 0.1 & 0.1 & 0.8 \end{bmatrix}$$

$$2) \mathcal{U}_0 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

$$\therefore P^\pi \mathcal{U}_0 = (I - \gamma P^\pi)^{-1} \mathcal{U}_0$$

```

1 import numpy as np
2
3 # Define the discount factor, initial state distribution, and probability transition matrix
4 gamma = 0.95
5 mu_0 = np.array([[1], [0], [0]])
6 P_pi = np.array([[0.1, 0.8, 0.1],
7                  [0.1, 0.1, 0.8],
8                  [0.1, 0.1, 0.8]])
9
10 # Compute the transpose of P_pi
11 P_pi_t = P_pi.T
12
13 # Compute (I - gamma * P_pi^T)
14 I = np.eye(3)
15 I_minus_gamma_P_pi_t = I - gamma * P_pi_t
16
17 # Compute (I - gamma * P_pi^T)^-1
18 I_minus_gamma_P_pi_t_inv = np.linalg.inv(I_minus_gamma_P_pi_t)
19
20 # Compute rho_pi_mu_0 = (I - gamma * P_pi^T)^-1 * mu_0
21 rho_pi_mu_0 = np.dot(I_minus_gamma_P_pi_t_inv, mu_0)
22
23 print("State occupancy measure:")
24 for i, rho in enumerate(rho_pi_mu_0, 1):
25     print(f"rho pi mu_0 ({i}) = {rho[0]:.4f}")

```

$$\therefore P^\pi \mathcal{U}_0 = \begin{bmatrix} 2.9 \\ 3.829 \\ 13.272 \end{bmatrix}$$

$$3) \bar{P}^{\pi} N_0 = \frac{P^{\pi} N_0}{\sum P^{\pi} N_0} = \begin{bmatrix} 0.1450 \\ 0.1914 \\ 0.6636 \end{bmatrix}$$

4) To find $\|P'(\cdot | s, \pi(s)) - P(\cdot | s, \pi(s))\|_1$,

```

1 import numpy as np
2
3 # original transition probabilities under policy pi
4 P_pi = np.array([[0.1, 0.8, 0.1], [0.1, 0.1, 0.8], [0.1, 0.1, 0.8]])
5
6 # new transition probabilities under policy pi
7 P_prime_pi = np.array([[0.15, 0.7, 0.15], [0.05, 0.05, 0.9], [0.15, 0.8, 0.05]])
8
9 # compute the L1-difference for each state
10 l1_diff = np.sum(np.abs(P_prime_pi - P_pi), axis=1)
11
12 # print the L1-difference for each state
13 for i, diff in enumerate(l1_diff, 1):
14     print(f"s = {i}: |P'(.|{i}, pi({i})) - P(.|{i}, pi({i}))|_1 = {diff:.1f}")
15

```

$$\left. \begin{aligned} \therefore s \leftarrow 1 &= 0.2 \\ s \leftarrow 2 &= 0.2 \\ s \leftarrow 3 &= 1.5 \end{aligned} \right\} L_s$$

5) For $s=1$, find $|V'^\pi(1) - V^\pi(1)|$

We can use the previously calculated $P^\pi V_0$ & $\bar{P}^\pi V_0$

$$: |V'^\pi(1) - V^\pi(1)| = \frac{1}{1-\gamma} \sum \bar{P}^\pi V_0 \cdot L_s$$

\downarrow
 $L_1 \text{ diff}$

```
1 import numpy as np
2
3 # original transition probabilities under policy pi
4 P_pi = np.array([[0.1, 0.8, 0.1], [0.1, 0.1, 0.8], [0.1, 0.1, 0.8]])
5
6 # new transition probabilities under policy pi
7 P_prime_pi = np.array([[0.15, 0.7, 0.15], [0.05, 0.05, 0.9], [0.15, 0.8, 0.05]])
8
9 # compute the L1-difference for each state
10 l1_diff = np.sum(np.abs(P_prime_pi - P_pi), axis=1)
11
12 # print the L1-difference for each state
13 for i, diff in enumerate(l1_diff, 1):
14     print(f"s = {i}: |P'(.|{i}, pi({i})) - P(.|{i}, pi({i}))|_1 = {diff:.1f}")
15
```

$$\therefore |V'^\pi(1) - V^\pi(1)| = \underline{\underline{21.2536}}$$

Problem 2:

$$1) R(s, a) = \begin{cases} 1 & , (3, h) \\ 0 & \end{cases}$$

$$\therefore \pi(1) = g, \pi(2) = g, \pi(3) = h$$

$$\therefore P = \begin{bmatrix} 0.8 & 0.1 & 0.1 \\ 0.1 & 0.8 & 0.1 \\ 0.1 & 0.1 & 0.8 \end{bmatrix}$$

$$R = [[0, 0], [0, 0], [0, 1]]$$



```
1  import numpy as np
2
3  # Transition probabilities and rewards
4  P = np.array(
5      [
6          [[0.1, 0.8, 0.1], [0.8, 0.1, 0.1]],
7          [[0.1, 0.1, 0.8], [0.1, 0.8, 0.1]],
8          [[0.8, 0.1, 0.1], [0.1, 0.1, 0.8]],
9      ]
10 )
11 R = np.array([[0, 0], [0, 0], [0, 1]])
12
13 # Policy
14 pi = np.array([0, 0, 1])
15
16 gamma = 0.95
17
18 # set up system of linear equations
19 A = np.eye(3) - gamma * P[np.arange(3), pi]
20 b = R[np.arange(3), pi]
21
22 # solve system of linear equations
23 V_pi = np.linalg.solve(A, b)
24
25 for i, v in enumerate(V_pi, 1):
26     print(f"V  $\pi(\{i\}) \approx {v:.4f}")$ 
27
```

2) Given oracle access to $P(s'|s,a)$, we can generate 100 samples of s' using each s,a pair, count the number of times each s' appears, & divide by 100 to obtain $\hat{P}(s'|s,a)$

& Assuming oracle access in code \rightarrow using given transition probabilities to sample

```
1 import numpy as np
2
3 P_true = np.array([
4     [[0.1, 0.8, 0.1], [0.8, 0.1, 0.1]],
5     [[0.1, 0.1, 0.8], [0.1, 0.8, 0.1]],
6     [[0.8, 0.1, 0.1], [0.1, 0.1, 0.8]]
7 ])
8
9 num_samples = 100
10 states = [1, 2, 3]
11 actions = ['g', 'h']
12
13 # initialize the estimated transition function
14 P_hat = np.zeros((3, 2, 3))
15
16 # sample each state-action pair and estimate the transition probabilities
17 for s in states:
18     for a_idx, a in enumerate(actions):
19         # generate samples using the true transition probabilities
20         samples = np.random.choice([1, 2, 3], size=num_samples, p=P_true[s - 1, a_idx])
21
22         # count the occurrences of each next state in the samples
23         for s_prime in states:
24             count = np.sum(samples == s_prime)
25             P_hat[s - 1, a_idx, s_prime - 1] = count / num_samples
26
27 print("Estimated transition function P^:")
28 for s in states:
29     for a_idx, a in enumerate(actions):
30         print(f"P^(.{s}, {a}) = {P_hat[s - 1, a_idx]}")
```


● PS C:\Users\suddu\Downloads\Code\ECE\495_RL\hw3> & C:/Users/suddu/AppData/Local/Programs/Python/Python312/python.exe c:/Users/suddu/Downloads/Code/ECE/495_RL/hw3/hw3_2.2.py

Estimated transition function P^{\cdot} :

$P^{\cdot}(.|1, g) = [0.09 \ 0.84 \ 0.07]$

$P^{\cdot}(.|1, h) = [0.85 \ 0.09 \ 0.06]$

$P^{\cdot}(.|2, g) = [0.11 \ 0.05 \ 0.84]$

$P^{\cdot}(.|2, h) = [0.11 \ 0.79 \ 0.1 \]$

$P^{\cdot}(.|3, g) = [0.75 \ 0.06 \ 0.19]$

$P^{\cdot}(.|3, h) = [0.13 \ 0.14 \ 0.73]$

○ PS C:\Users\suddu\Downloads\Code\ECE\495_RL\hw3> █

Question 2.3

```
1 import numpy as np
2
3 P_true = np.array([
4     [[0.1, 0.8, 0.1], [0.8, 0.1, 0.1]],
5     [[0.1, 0.1, 0.8], [0.1, 0.8, 0.1]],
6     [[0.8, 0.1, 0.1], [0.1, 0.1, 0.8]]
7 ])
8
9 num_samples = 100
10 states = [1, 2, 3]
11 actions = ['g', 'h']
12
13 # initialize the estimated transition function
14 P_hat = np.zeros((3, 2, 3))
15
16 # sample each s, a pair and estimate the transition probabilities
17 for s in states:
18     for a_idx, a in enumerate(actions):
19         # generate samples using the true transition probabilities
20         samples = np.random.choice([1, 2, 3], size=num_samples, p=P_true[s - 1, a_idx])
21
22         # count the occurrences of each s' in the samples
23         for s_prime in states:
24             count = np.sum(samples == s_prime)
25             P_hat[s - 1, a_idx, s_prime - 1] = count / num_samples
26
27 print("Estimated transition function P^:")
28 for s in states:
29     for a_idx, a in enumerate(actions):
30         print(f"P^{s}, {a} = {P_hat[s - 1, a_idx]}")
31
32 pi = ['g', 'g', 'h']
33
34 # compute the L1-difference for each state under policy pi
35 l1_diff = np.zeros(3)
36 for s in range(3):
37     a_idx = 0 if pi[s] == 'g' else 1
38     l1_diff[s] = np.sum(np.abs(P_hat[s, a_idx] - P_true[s, a_idx]))
39
40 print("\nL1-difference under policy pi:")
41 for s in range(3):
42     print(f"|P^{s+1}, {pi[s]} - P^{s+1}, {pi[s]}|_1 = {l1_diff[s]:.4f}")
```

$$\therefore \|P^{\wedge}(\cdot | 1, g) - P(\cdot | 1, g)\|_1 = \underline{\underline{0.08}}$$

$$\|P^{\wedge}(\cdot | 2, g) - P(\cdot | 2, g)\|_1 = \underline{\underline{0.1}}$$

$$\|P^{\wedge}(\cdot | 3, h) - P(\cdot | 3, h)\|_1 = \underline{\underline{0.1}}$$

Question 2.5

```
1 import numpy as np
2
3 P_true = np.array([
4     [[0.1, 0.8, 0.1], [0.8, 0.1, 0.1]],
5     [[0.1, 0.1, 0.8], [0.1, 0.8, 0.1]],
6     [[0.8, 0.1, 0.1], [0.1, 0.1, 0.8]]
7 ])
8
9 num_samples = 100
10 states = [1, 2, 3]
11 actions = ['g', 'h']
12
13 # initialize the estimated transition function
14 P_hat = np.zeros((3, 2, 3))
15
16 # sample each s, a pair and estimate the transition probabilities
17 for s in states:
18     for a_idx, a in enumerate(actions):
19         # generate samples using the true transition probabilities
20         samples = np.random.choice([1, 2, 3], size=num_samples, p=P_true[s - 1, a_idx])
21
22         # count the occurrences of each s' in the samples
23         for s_prime in states:
24             count = np.sum(samples == s_prime)
25             P_hat[s - 1, a_idx, s_prime - 1] = count / num_samples
26
27 print("Estimated transition function P̂:")
28 for s in states:
29     for a_idx, a in enumerate(actions):
30         print(f"P^(.|{s}, {a}) = {P_hat[s - 1, a_idx]}")
31
32 pi = ['g', 'g', 'h']
33
34 # compute the L1-difference for each state under policy pi
35 l1_diff = np.zeros(3)
36 for s in range(3):
37     a_idx = 0 if pi[s] == 'g' else 1
38     l1_diff[s] = np.sum(np.abs(P_hat[s, a_idx] - P_true[s, a_idx]))
39
40 print("\nL1-difference under policy pi:")
41 for s in range(3):
42     print(f"P^(.|{s+1}, {pi[s]}) - P^(.|{s+1}, {pi[s]})|_1 = {l1_diff[s]:.4f}")
43
44 R = np.array([[0, 0], [0, 0], [0, 1]])
45
46 gamma = 0.95
47
48 # set up the system of linear equations
49 A = np.eye(3)
50 b = np.zeros(3)
51
52 for s in range(3):
53     a_idx = 0 if pi[s] == 'g' else 1
54     A[s] -= gamma * P_hat[s, a_idx]
55     b[s] = R[s, a_idx]
56
57 # solve system of linear equations
58 V_hat_pi = np.linalg.solve(A, b)
59
60 print("\nEstimated value function V^ pi:")
61 for s in range(3):
62     print(f"V^ pi({s+1}) = {V_hat_pi[s]:.4f}")
```

$$\therefore \hat{V}^{\pi}(1) = 12.955$$

$$\hat{V}^{\pi}(2) = 13.650$$

$$\hat{V}^{\pi}(3) = 15.395 //$$

$$6) |\hat{V}^{\pi}(1) - V^{\pi}(1)| = \underline{\underline{12.8841}}$$

4) Bonus

We can use the following formula

$$|\hat{V}^{\pi}(1) - V^{\pi}(1)| \leq \frac{1}{1-\gamma} \sum \bar{\rho}^{\pi} \mu_0(s) \|\hat{P}(\cdot|s, \pi(s)) - P(\cdot|s, \pi(s))\|,$$

```

1 rho_bar_pi_mu_0 = np.array([2.9, 3.8285, 13.2715])
2
3 # compute the upper bound using the simulation lemma
4 upper_bound = (1 / (1 - gamma)) * np.sum(rho_bar_pi_mu_0 * l1_diff)
5
6 print(f"\nUpper bound on |V^ pi(1) - V pi(1)| using the simulation lemma:")
7 print(f"|V^ pi(1) - V pi(1)| <= {upper_bound:.4f}")
8

```

$$\therefore \|\hat{V}^{\pi}(1) - V^{\pi}(1)\|_1 \leq \underline{\underline{13.3828}}$$

Problem 3:

1) To show that the tabular representations of $Q(s, a)$ is a special class of linear functions, we need to define the number of features k & the vector val. function $\phi(s, a)$ s.t the lin. func. is eq. to the tabular representation.

$$\therefore k = |S| \times |A|$$

$$\phi(s, a) : S \times A \rightarrow \mathbb{R}^k$$

$$\therefore Q = \{Q_\theta : Q_\theta(s, a) = \theta^T \phi(s, a) \text{ for all } s \in S \text{ \& } a \in A, \theta \in \mathbb{R}^k\}$$

For any (s, a) , $\phi(s, a)$ is a vector with a 1 in the position corresponding to (s, a) & 0's elsewhere.

$$\therefore Q_\theta(s, a) = \theta^T \phi(s, a) = \underbrace{\theta^*(s, a)}$$

* elements of parameter vector θ corr. to (s, a)

This shows that the linear func. approx. $Q_0(s,a)$ with the defined feature function $\phi(s,a)$ & number of features $k = |S| \times |A|$ is eq. to the tabular representation, where (s,a) pairs have their own indep. parameters $\theta^*(s,a)$

\therefore The # of features k should be equal to the total number (s,a) pairs

\rightarrow The feature function $\phi(s,a)$ should map each s,a pair to a binary vector with only one 1 of size k , positioned corresponding to (s,a) .

With these definitions, the resulting class of linear function Q parameterized by θ is equivalent to the tabular representations of the s,a value function.

2) $k = n \times m$ where n : # of state partitions
 m : # of action partitions

$\phi(s,a) = e\{j,i\}$ is a one-hot vector of size k with a 1 in the position corresponding to the state partition index j & action partition index i , & 0's elsewhere.

$\therefore s \in S_j$ & $a \in A_i, \phi(s,a) = e\{j,i\}$

$\therefore Q_\theta(s,a) = \theta^T \phi(s,a)$ is the lin. func. approx.

$$= \theta^T e\{j,i\} = \underbrace{\theta\{j,i\}}$$

element of the parameter vector θ corresponding to the state partition index j & action index i

This func. shows that the linear func. approx. $Q_\theta(s,a)$ with the defined feature function $\phi(s,a)$ & # of features $R = n \times m$ assigns the same val. $\theta \{i, j\}$ to all s, a pairs.

This representation is more efficient than the tabular one because it reduces the parameter size from $|S| \times |A|$ to $n \times m$
of state & action partitions

By grouping similar states & actions into partitions we can capture the shared structure in Q func. & reduce its parameter size. //

3) Bonus

Since we have a linear func. approx. with $\phi(s,a)$ parameterized by θ :

$$Q^n(s,a) = \operatorname{argmin}_{Q \in \mathcal{Q}} \sum_{i=1}^N (y_i - Q(s_i, a_i))^2$$

where \mathcal{Q} represents the class of linear functions based on the single feature $\phi(s, a)$

$$Q(s, a) = \theta \phi(s, a) \rightarrow \text{lin. func. approx.}$$

$$\therefore Q^n(s, a) = \arg \min_{\theta} \sum_{i=1}^N (y_i^2 - 2y_i \theta \phi(s_i, a_i) + \theta^2 \phi(s_i, a_i)^2)$$

Since $\phi(s_i, a_i)$ is known & indep. of θ , we can treat it as constant

$$\therefore \frac{\partial \mathcal{Q}}{\partial \theta} = 0$$

$$\therefore \sum_{i=1}^N (-2y_i \phi(s_i, a_i) + 2\theta \phi(s_i, a_i)^2) = 0$$

$$\rightarrow -2 \sum_{i=1}^N y_i \phi(s_i, a_i) + 2\theta \sum_{i=1}^N \phi(s_i, a_i)^2 = 0$$

$$\rightarrow \theta \sum_{i=1}^N \phi(s_i, a_i)^2 = \sum_{i=1}^N y_i \phi(s_i, a_i)$$

$$\therefore \theta = \frac{\sum_{i=1}^N y_i \phi(s_i, a_i)}{\sum_{i=1}^N \phi(s_i, a_i)^2}$$

$$\sum_{i=1}^N \phi(s_i, a_i)^2 //$$