



## Problem 1.

Consider a system with  $\gamma=0.9$  the following state and action spaces:  
 $S = \{-1, 1, 2\}$ ,  $A = \{-1, 0, 1\}$ . The available batch data are as follows:

$$D = \{ (s_0=1, a_0=1, r_1=1, s_1=2), (s_1=2, a_1=0, r_2=-1, s_2=1), \\ (s_2=1, a_2=-1, r_3=0, s_3=-1) \}$$

consider the basis function  $\Phi(s, a) = a^2 s + a s \rightarrow a$  with initial weights  $w^0 = 1$ .  
 Perform LSPI to compute  $w^1$  and  $w^2$ , and policy associated to  $w^2$ .

\* In case of tie for action selection, give the preference to -1, then 0 and finally 1.

Example  $\rightarrow \arg \max_{a \in \{-1, 0, 1\}} \left\{ \frac{-1}{2} \quad \frac{0}{2} \quad \frac{+1}{2} \right\} = -1$

### Problem 2.

Repeat Problem 1 using the basis function  $\Phi(s, a) = \begin{bmatrix} a s + a \\ a^2 s \end{bmatrix}$  with initial weights  $w^0 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$ . Perform LSPI to compute  $w^1$  and  $w^2$ , and policy associated to  $w^2$ . Is the final policy (i.e.,  $\pi^2$ ) different from Problem 1? Can two basis functions, in general, lead to different policies?

### Problem 3.

Consider a system with the following continuous state and action spaces:  $S = [-2, 2]$ ,  $A = [-1.5, 2]$ . The available batch data is as:

$$D = \{ (s_0 = 2, a_0 = 1, r_1 = 1, s_1 = 1), (s_1 = 1, a_1 = -1, r_2 = 2, s_2 = -1) \}$$

Consider  $\gamma = 0.9$ , the basis function  $\Phi(s, a) = \begin{bmatrix} a \\ s \times a \end{bmatrix}$  with initial weights  $w^0 = \begin{bmatrix} 0.5 \\ 1 \end{bmatrix}$ . Perform LSPI to compute  $w^1$  and  $w^2$ . Compute the policy associated with  $w^2$  for any  $s \in [-1, 1]$ .

# Least Squares Policy Iteration

Batch Data:  $D = \{(s_1, a_1, r_2, s_2), \dots, (s_L, a_L, r_{L+1}, s_{L+1})\}$

Policy Evaluation

$$\underset{k \times 1}{\omega^-} \rightarrow \underset{k \times k}{Q}(s, a) = \underset{k \times 1}{\Phi}(s, a)^T \underset{k \times 1}{\omega^-}$$

$$\underset{k \times 1}{\pi^-}(s) = \underset{a \in A}{\operatorname{argmax}} Q(s, a) = \underset{a \in A}{\operatorname{argmax}} \underset{k \times 1}{\Phi}(s, a)^T \underset{k \times 1}{\omega^-}$$

$$\underset{k \times 1}{\omega^-} = \underset{k \times 1}{\omega^+}$$

Policy Improvement

$$\underset{k \times k}{A} = \frac{1}{L} \sum_{i=1}^L \underset{k \times 1}{\Phi}(s_i, a_i) \left[ \underset{k \times 1}{\Phi}(s_i, a_i) - \gamma \underset{k \times 1}{\Phi}(s_{i+1}, \pi^-(s_{i+1})) \right]^T$$

$$\underset{k \times 1}{b} = \frac{1}{L} \sum_{i=1}^L \underset{k \times 1}{\Phi}(s_i, a_i) r_{i+1}$$

$$\underset{k \times 1}{\omega^+} = \underset{k \times k}{A}^{-1} \underset{k \times 1}{b}$$

Questions about the HW should be directed to TA, Begum Taskazan, at [taskazan.b@northeastern.edu](mailto:taskazan.b@northeastern.edu).