

Homework Set 4

Introduction to Artificial Intelligence with Mathematics (MAS473)

Total Points = 50pts

1. (15pts) In this problem, we consider a toy example of applying a Markov Chain Monte Carlo (MCMC) method. Let X be a discrete random variable whose unnormalized probability mass function is

$$\tilde{p}(X = 1) = 4, \tilde{p}(X = 2) = 2, \tilde{p}(X = 3) = 3, \tilde{p}(X = 4) = 3.$$

Our goal is to evaluate $\mathbf{E}[f(X)]$ where $f(x) = x^2 - 3$. To apply the Metropolis-Hasting algorithm, define two proposal distributions $q(i, j) (= q_{ij})$ and $\tilde{q}(i, j) (= \tilde{q}_{ij})$ where

$$q(i, i+1) = q(i, i-1) = \frac{1}{2}, \quad \tilde{q}(i, i+1) = \frac{2}{3}, \tilde{q}(i, i-1) = \frac{1}{3}.$$

(If $i+1 = 5$, consider 5 as 1. Also, if $i-1 = 0$, consider 0 as 4.)

- (a) Using the Metropolis-Hasting algorithm, construct Markov chains for proposal distributions q and \tilde{q} when we set the initial state $X_0 = 1$. You should find the transition matrices for these Markov chains.
- (b) For q and \tilde{q} , make a python code to evaluate the error between the actual value of $\mathbf{E}[f(X)]$ and $\mathbf{E}\left[\frac{1}{n} \sum_{i=0}^{n-1} f(X_i)\right]$ where $\{X_i\}$ is the Markov chain from (a) and plot the errors along $n = 1, \dots, 1000$. Which proposal distribution converges faster in the expected sense? You should calculate the exact $\mathbf{E}[f(X)]$ and $\mathbf{E}\left[\frac{1}{n} \sum_{i=0}^{n-1} f(X_i)\right]$, not a sample mean of $f(X)$ and $\frac{1}{n} \sum_{i=0}^{n-1} f(X_i)$.

Solution.

- (a) The transition matrix $P(X_{t+1} = j | X_t = i)$ of the Markov chain X_t when we use the proposal distribution q is

	$j = 1$	$j = 2$	$j = 3$	$j = 4$
$i = 1$	$1 - \frac{1}{4} - \frac{3}{8} = \frac{3}{8}$	$\frac{1}{2} \times \min\left(\frac{1/6 \times 1/2}{1/3 \times 1/2}, 1\right) = \frac{1}{4}$	0	$\frac{1}{2} \times \min\left(\frac{1/4 \times 1/2}{1/3 \times 1/2}, 1\right) = \frac{3}{8}$
$i = 2$	$\frac{1}{2} \times \min\left(\frac{1/3 \times 1/2}{1/6 \times 1/2}, 1\right) = \frac{1}{2}$	$1 - \frac{1}{2} - \frac{1}{2} = 0$	$\frac{1}{2} \times \min\left(\frac{1/4 \times 1/2}{1/6 \times 1/2}, 1\right) = \frac{1}{2}$	0
$i = 3$	0	$\frac{1}{2} \times \min\left(\frac{1/6 \times 1/2}{1/4 \times 1/2}, 1\right) = \frac{1}{3}$	$1 - \frac{1}{3} - \frac{1}{2} = \frac{1}{6}$	$\frac{1}{2} \times \min\left(\frac{1/4 \times 1/2}{1/4 \times 1/2}, 1\right) = \frac{1}{2}$
$i = 4$	$\frac{1}{2} \times \min\left(\frac{1/3 \times 1/2}{1/4 \times 1/2}, 1\right) = \frac{1}{2}$	0	$\frac{1}{2} \times \min\left(\frac{1/4 \times 1/2}{1/4 \times 1/2}, 1\right) = \frac{1}{2}$	$1 - \frac{1}{2} - \frac{1}{2} = 0$

The transition matrix $P(X_{t+1} = j | X_t = i)$ of the Markov chain X_t when we use the proposal distribution \tilde{q} is

	$j = 1$	$j = 2$	$j = 3$	$j = 4$
$i = 1$	$1 - \frac{1}{6} - \frac{1}{3} = \frac{1}{2}$	$\frac{2}{3} \times \min\left(\frac{1/6 \times 1/3}{1/3 \times 2/3}, 1\right) = \frac{1}{6}$	0	$\frac{1}{3} \times \min\left(\frac{1/4 \times 2/3}{1/3 \times 1/3}, 1\right) = \frac{1}{3}$
$i = 2$	$\frac{1}{3} \times \min\left(\frac{1/3 \times 2/3}{1/6 \times 1/3}, 1\right) = \frac{1}{3}$	$1 - \frac{1}{3} - \frac{1}{2} = \frac{1}{6}$	$\frac{2}{3} \times \min\left(\frac{1/4 \times 1/3}{1/6 \times 2/3}, 1\right) = \frac{1}{2}$	0
$i = 3$	0	$\frac{1}{3} \times \min\left(\frac{1/6 \times 2/3}{1/4 \times 1/3}, 1\right) = \frac{1}{3}$	$1 - \frac{1}{3} - \frac{1}{3} = \frac{1}{3}$	$\frac{2}{3} \times \min\left(\frac{1/4 \times 1/3}{1/4 \times 2/3}, 1\right) = \frac{1}{3}$
$i = 4$	$\frac{2}{3} \times \min\left(\frac{1/3 \times 1/3}{1/4 \times 2/3}, 1\right) = \frac{4}{9}$	0	$\frac{1}{3} \times \min\left(\frac{1/4 \times 2/3}{1/4 \times 1/3}, 1\right) = \frac{1}{3}$	$1 - \frac{4}{9} - \frac{1}{3} = \frac{2}{9}$

(b) First,

$$\mathbf{E}[f(X)] = \frac{4}{12} \times (1^2 - 3) + \frac{2}{12} \times (2^2 - 3) + \frac{3}{12} \times (3^2 - 3) + \frac{3}{12} \times (4^2 - 3) = \frac{51}{12}.$$

Note that

$$\mathbf{E} \left[\frac{1}{n} \sum_{i=0}^{n-1} f(X_i) \right] = \frac{1}{n} \sum_{i=0}^{n-1} \mathbf{E}[f(X_i)]$$

and the distribution of X_i is

$$\begin{bmatrix} 1 & 0 & 0 & 0 \end{bmatrix} P^i$$

where P is the transition matrix of the Markov chain X_t . Using the code implementation, the Markov chain from the Metropolis-Hasting algorithm with q is slightly faster than another one with \tilde{q} . (See the .ipynb solution file.)

2. (15pts) Consider a Markov decision process (S, A, P, R) where $S = \{0, 1, 2, 3, 4\}$ is a state space, $A = \{a^1, a^2\}$ is an action space, P is a transition probability matrix such that

$$P((s+1)\%5|s, a^1) = P((s+2)\%5|s, a^1) = P((s+4)\%5|s, a^1) = \frac{1}{3},$$

$$P((s+1)\%5|s, a^2) = P((s+3)\%5|s, a^2) = \frac{1}{2}$$

(where $a\%b$ means the remainder when a is divided by b) for all $s \in S$ and R is the reward such that $R(s, a, s')$ follows Bernoulli($\frac{1}{2}$) if $s' = 0$ and 0 otherwise.¹ Also, when the process reaches the state 4, then the process is terminated (i.e. 4 is the terminal state). Assume the initial state is $s_0 = 0$ and set the discounted factor $\gamma = 0.9$. Let π be a Markovian randomized stationary policy that $\pi(a^1|s) = \pi(a^2|s) = 0.5$ for $s = 0, 2$ and $\pi(a^1|s) = 0.7, \pi(a^2|s) = 0.3$ for $s = 1, 3$. In this problem, you may use the python code for matrix calculations, e.g. matrix addition, multiplication, inversion, ...

- (a) When we adopt the policy π , find the probability that the following trajectory is sampled:

$$\tau = (s_0 = 0, a_0 = a^1, r_0 = 0, s_1 = 2, a_1 = a^2, r_1 = 1, s_2 = 0, a_2 = a^2, r_2 = 0, s_3 = 3, a_3 = a^2, r_3 = 0, s_4 = 4)$$

- (b) Calculate $V^\pi(s)$ and $Q^\pi(s, a)$.
- (c) Let d_0 be a Markovian deterministic stationary policy that $d_0(s) = a^1$ for all $s \in S$. Using the policy iteration algorithm three times with the initial policy d_0 , show your result, e.g. policy evaluation result $v^{(n)}$ and the improved policy d_{n+1} .
- (d) Find all Markovian deterministic optimal policies and the optimal value function $V^*(s)$.

Solution.

- (a)

$$\begin{aligned} p(\tau) &= p(s_0 = 0) \cdot \pi(a_0 = a^1|s_0 = 0) \cdot p(s_1 = 2|s_0 = 0, a_0 = a^1) \cdot p(r_0 = 0|s_0 = 0, a_0 = a^1, s_1 = 2) \\ &\quad \cdot \pi(a_1 = a^2|s_1 = 2) \cdot p(s_2 = 0|s_1 = 2, a_1 = a^2) \cdot p(r_1 = 1|s_1 = 2, a_1 = a^2, s_2 = 0) \\ &\quad \cdot \pi(a_2 = a^2|s_2 = 0) \cdot p(s_3 = 3|s_2 = 0, a_2 = a^2) \cdot p(r_2 = 0|s_2 = 0, a_2 = a^2, s_3 = 3) \\ &\quad \cdot \pi(a_3 = a^2|s_3 = 3) \cdot p(s_4 = 4|s_3 = 3, a_3 = a^2) \cdot p(r_3 = 0|s_3 = 3, a_3 = a^2, s_4 = 4) \\ &= 1 \times \frac{1}{2} \times \frac{1}{3} \times 1 \times \frac{1}{2} \times \frac{1}{2} \times \frac{1}{2} \times \frac{1}{2} \times 1 \times \frac{3}{10} \times \frac{1}{2} \times 1 \\ &= \frac{1}{1280} \end{aligned}$$

- (b) Since $r(s, a, s')$ is independent of s and a ,

$$V^\pi(s) = \mathbf{E}_{a \sim \pi(\cdot|s), s' \sim p(\cdot|s, a)} [r(s, a, s') + \gamma V^\pi(s')] = \mathbf{E}_{a \sim \pi(\cdot|s), s' \sim p(\cdot|s, a)} [r(s') + \gamma V^\pi(s')].$$

Under the policy π , $\{s_t\}$ is a Markov chain with a transition matrix $p(s_{t+1} = j|s_t = i)$ is

¹In the regular lecture, we use $R(s, a)$ but we can consider a general reward $R(s, a, s')$. It means that the reward when the state is s , the action a is taken and s' is the next state.

	$j = 0$	$j = 1$	$j = 2$	$j = 3$	$j = 4$
$i = 0$	0	$\frac{1}{2} \times \frac{1}{3} + \frac{1}{2} \times \frac{1}{2} = \frac{5}{12}$	$\frac{1}{2} \times \frac{1}{3} = \frac{1}{6}$	$\frac{1}{2} \times \frac{1}{2} = \frac{1}{4}$	$\frac{1}{2} \times \frac{1}{3} = \frac{1}{6}$
$i = 1$	$\frac{7}{10} \times \frac{1}{3} = \frac{7}{30}$	0	$\frac{7}{10} \times \frac{1}{3} + \frac{3}{10} \times \frac{1}{2} = \frac{23}{60}$	$\frac{7}{10} \times \frac{1}{3} = \frac{7}{30}$	$\frac{3}{10} \times \frac{1}{2} = \frac{3}{20}$
$i = 2$	$\frac{1}{2} \times \frac{1}{2} = \frac{1}{4}$	$\frac{1}{2} \times \frac{1}{3} = \frac{1}{6}$	0	$\frac{1}{2} \times \frac{1}{3} + \frac{1}{2} \times \frac{1}{2} = \frac{5}{12}$	$\frac{1}{2} \times \frac{1}{3} = \frac{1}{6}$
$i = 3$	$\frac{7}{10} \times \frac{1}{3} = \frac{7}{30}$	$\frac{3}{10} \times \frac{1}{2} = \frac{3}{20}$	$\frac{7}{10} \times \frac{1}{3} = \frac{7}{30}$	0	$\frac{7}{10} \times \frac{1}{3} + \frac{3}{10} \times \frac{1}{2} = \frac{23}{60}$
$i = 4$	0	0	0	0	1

Then the Bellman equation is

$$\begin{bmatrix} V^\pi(0) \\ V^\pi(1) \\ V^\pi(2) \\ V^\pi(3) \\ V^\pi(4) \end{bmatrix} = \begin{bmatrix} 0 \\ \frac{7}{60} \\ \frac{1}{8} \\ \frac{7}{60} \\ 0 \end{bmatrix} + 0.9 \cdot \begin{bmatrix} 0 & \frac{5}{12} & \frac{1}{6} & \frac{1}{4} & \frac{1}{6} \\ \frac{7}{30} & 0 & \frac{23}{60} & \frac{7}{30} & \frac{3}{20} \\ \frac{1}{4} & \frac{1}{6} & 0 & \frac{5}{12} & \frac{1}{6} \\ \frac{7}{30} & \frac{3}{20} & \frac{7}{30} & 0 & \frac{23}{60} \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} V^\pi(0) \\ V^\pi(1) \\ V^\pi(2) \\ V^\pi(3) \\ V^\pi(4) \end{bmatrix}$$

Using a python library, we can conclude

$$\begin{bmatrix} V^\pi(0) \\ V^\pi(1) \\ V^\pi(2) \\ V^\pi(3) \\ V^\pi(4) \end{bmatrix} = \begin{bmatrix} 0.24447704 \\ 0.34493336 \\ 0.33890058 \\ 0.28574197 \\ 0 \end{bmatrix}.$$

Also, we observe that

$$Q^\pi(s, a) = \mathbf{E}_{s' \sim p(\cdot | s, a)} [r(s, a, s') + \gamma V^\pi(s')].$$

Then,

$$\begin{bmatrix} Q^\pi(0, a^1) \\ Q^\pi(1, a^1) \\ Q^\pi(2, a^1) \\ Q^\pi(3, a^1) \end{bmatrix} = \begin{bmatrix} 0 \\ \frac{1}{6} \\ 0 \\ \frac{1}{6} \end{bmatrix} + 0.9 \cdot \begin{bmatrix} 0 & \frac{1}{3} & \frac{1}{3} & 0 & \frac{1}{3} \\ \frac{1}{3} & 0 & \frac{1}{3} & \frac{1}{3} & 0 \\ 0 & \frac{1}{3} & 0 & \frac{1}{3} & \frac{1}{3} \\ \frac{1}{3} & 0 & \frac{1}{3} & 0 & \frac{1}{3} \end{bmatrix} \begin{bmatrix} V^\pi(0) \\ V^\pi(1) \\ V^\pi(2) \\ V^\pi(3) \\ V^\pi(4) \end{bmatrix} = \begin{bmatrix} 0.20515018 \\ 0.42740254 \\ 0.1892026 \\ 0.34167995 \end{bmatrix}.$$

Also,

$$\begin{bmatrix} Q^\pi(0, a^2) \\ Q^\pi(1, a^2) \\ Q^\pi(2, a^2) \\ Q^\pi(3, a^2) \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \frac{1}{4} \\ 0 \end{bmatrix} + 0.9 \cdot \begin{bmatrix} 0 & \frac{1}{2} & 0 & \frac{1}{2} & 0 \\ 0 & 0 & \frac{1}{2} & 0 & \frac{1}{2} \\ \frac{1}{2} & 0 & 0 & \frac{1}{2} & 0 \\ 0 & \frac{1}{2} & 0 & 0 & \frac{1}{2} \end{bmatrix} \begin{bmatrix} V^\pi(0) \\ V^\pi(1) \\ V^\pi(2) \\ V^\pi(3) \\ V^\pi(4) \end{bmatrix} = \begin{bmatrix} 0.2838039 \\ 0.15250526 \\ 0.48859855 \\ 0.15522001 \end{bmatrix}.$$

(c) First, we calculate $v^{(0)}$ by the formula $v^{(0)} = (I - \gamma P_{d_0})^{-1} r_{d_0}$:

$$\begin{bmatrix} v^{(0)}(0) \\ v^{(0)}(1) \\ v^{(0)}(2) \\ v^{(0)}(3) \\ v^{(0)}(4) \end{bmatrix} = \left(I - 0.9 \begin{bmatrix} 0 & \frac{1}{3} & \frac{1}{3} & 0 & \frac{1}{3} \\ \frac{1}{3} & 0 & \frac{1}{3} & \frac{1}{3} & 0 \\ 0 & \frac{1}{3} & 0 & \frac{1}{3} & \frac{1}{3} \\ \frac{1}{3} & 0 & \frac{1}{3} & 0 & \frac{1}{3} \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix} \right)^{-1} \begin{bmatrix} 0 \\ \frac{1}{6} \\ 0 \\ \frac{1}{6} \\ 0 \end{bmatrix} = \begin{bmatrix} 0.16207851 \\ 0.35293479 \\ 0.18732693 \\ 0.2714883 \\ 0 \end{bmatrix}$$

Then,

$$\begin{bmatrix} Q(0, a^1) \\ Q(1, a^1) \\ Q(2, a^1) \\ Q(3, a^1) \end{bmatrix} = \begin{bmatrix} 0.16207851 \\ 0.35293479 \\ 0.18732693 \\ 0.2714883 \end{bmatrix}.$$

If we evaluate $Q(s, a^2)$ using $v^{(0)}$,

$$\begin{bmatrix} Q(0, a^2) \\ Q(1, a^2) \\ Q(2, a^2) \\ Q(3, a^2) \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \frac{1}{4} \\ 0 \end{bmatrix} + 0.9 \cdot \begin{bmatrix} 0 & \frac{1}{2} & 0 & \frac{1}{2} & 0 \\ 0 & 0 & \frac{1}{2} & 0 & \frac{1}{2} \\ \frac{1}{2} & 0 & 0 & \frac{1}{2} & 0 \\ 0 & \frac{1}{2} & 0 & 0 & \frac{1}{2} \end{bmatrix} \begin{bmatrix} v^{(0)}(0) \\ v^{(0)}(1) \\ v^{(0)}(2) \\ v^{(0)}(3) \\ v^{(0)}(4) \end{bmatrix} = \begin{bmatrix} 0.28099039 \\ 0.08429712 \\ 0.44510507 \\ 0.15882065 \end{bmatrix}.$$

Then,

$$d^{(1)}(0) = a^2, d^{(1)}(1) = a^1, d^{(1)}(2) = a^2, d^{(1)}(3) = a^1.$$

Again, we calculate $v^{(1)}$:

$$\begin{bmatrix} v^{(1)}(0) \\ v^{(1)}(1) \\ v^{(1)}(2) \\ v^{(1)}(3) \\ v^{(1)}(4) \end{bmatrix} = \left(I - 0.9 \begin{bmatrix} 0 & \frac{1}{2} & 0 & \frac{1}{2} & 0 \\ \frac{1}{3} & 0 & \frac{1}{3} & \frac{1}{3} & 0 \\ \frac{1}{2} & 0 & 0 & \frac{1}{2} & 0 \\ \frac{1}{3} & 0 & \frac{1}{3} & 0 & \frac{1}{3} \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix} \right)^{-1} \begin{bmatrix} 0 \\ \frac{1}{6} \\ \frac{1}{4} \\ \frac{1}{6} \\ 0 \end{bmatrix} = \begin{bmatrix} 0.60303779 \\ 0.75743877 \\ 0.78355735 \\ 0.58264521 \\ 0 \end{bmatrix}$$

If we evaluate remainder $Q(s, a)$'s using $v^{(1)}$,

$$\begin{bmatrix} Q(0, a^1) \\ Q(1, a^2) \\ Q(2, a^1) \\ Q(3, a^2) \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} + 0.9 \cdot \begin{bmatrix} 0 & \frac{1}{3} & \frac{1}{3} & 0 & \frac{1}{3} \\ 0 & 0 & \frac{1}{2} & 0 & \frac{1}{2} \\ 0 & \frac{1}{3} & 0 & \frac{1}{3} & \frac{1}{3} \\ 0 & \frac{1}{2} & 0 & 0 & \frac{1}{2} \end{bmatrix} \begin{bmatrix} v^{(1)}(0) \\ v^{(1)}(1) \\ v^{(1)}(2) \\ v^{(1)}(3) \\ v^{(1)}(4) \end{bmatrix} = \begin{bmatrix} 0.46229884 \\ 0.35260081 \\ 0.40202519 \\ 0.34084745 \end{bmatrix}.$$

Thus, $d_2 = d_1$. Therefore, $d_3 = d_2 = d_1$ and $v^{(2)} = v^{(1)}$.

(d) We already found an optimal policy π^* in (c)

$$\pi^*(0) = a^2, \pi^*(1) = a^1, \pi^*(2) = a^2, \pi^*(3) = a^1.$$

Therefore, the optimal value function is

$$\begin{bmatrix} V^*(0) \\ V^*(1) \\ V^*(2) \\ V^*(3) \\ V^*(4) \end{bmatrix} = \begin{bmatrix} 0.60303779 \\ 0.75743877 \\ 0.78355735 \\ 0.58264521 \\ 0 \end{bmatrix}.$$

Suppose there exists another deterministic optimal policy $\tilde{\pi}$. Then

$$Q^{\tilde{\pi}}(s, \tilde{\pi}(s)) = \mathbf{E}_{s' \sim p(\cdot|s,a)} [r(s, a, s') + \gamma V^*(s')] = \mathbf{E}_{s' \sim p(\cdot|s,a)} [r(s, a, s') + \gamma V^{\pi^*}(s')] = Q^{\pi^*}(s, \tilde{\pi}(s)).$$

However, we already observed that $\pi^*(s)$ is the unique element of $\operatorname{argmin}_a Q^{\pi^*}(s, a)$ since

$$\begin{aligned} 0.60303779 &> 0.46229884, 0.75743877 > 0.35260081, \\ 0.78355735 &> 0.40202519, 0.58264521 > 0.34084745. \end{aligned}$$

Therefore, π^* is the unique Markovian deterministic optimal policy.

3. (5pts) In this problem, we solve a problem to find the shortest path from the Start to the End in the maze using the Markov decision process.

Start (1, 5)	(2, 5)			(5, 5)
(1, 4)	(2, 4)	(3, 4)	(4, 4)	(5, 4)
(1, 3)			(4, 3)	
(1, 2)	(2, 2)		(4, 2)	(5, 2)
	(2, 1)		(4, 1)	End (5, 1)

We formulate this problem as the following:

- The state space \mathcal{S} consists of the white blocks in the figure. Each element is expressed as (i, j) .
- The action space $A = \{N, S, W, E\}$ where N, S, W and E mean that go to north, south, west and east, respectively. However, the next state does not change if the direction of action is blocked. For example,

$$P((1, 3)|(1, 4), S) = 1, \quad P((4, 4)|(3, 4), E) = 1, \\ P((1, 5)|(1, 5), N) = 1, \quad P((4, 3)|(4, 4), S) = 1.$$

- $(1, 5)$ is the initial state and $(5, 1)$ is the terminal state.²
- $R(s, a, s') = 1$ if $s' = (5, 1)$ and 0 otherwise.
- The discounted factor $\gamma = 0.9$.

Initialize $V_0(s) = 0$ for all $s \in \mathcal{S}$. Using the Bellman optimal operator

$$(Lv)(s) = \sup_{a \in A} \mathbf{E}_{s' \sim P(\cdot|s,a)} [r(s, a, s') + \gamma \cdot v(s')],$$

set

$$V_{n+1}(s) = (LV_n)(s), \quad n = 0, 1, 2, \dots$$

A policy π is said to be a greedy policy induced by a value function V if

$$\pi(s) \in \operatorname{argmax}_{a \in A} Q(s, a)$$

where $Q(s, a) = \mathbf{E}_{s'} [R(s, a, s') + \gamma V(s')]$. Find the minimum $n \in \mathbb{N}$ such that all of greedy policies induced by V_n is an optimal policy. Note that π^* is said to be an optimal policy if $V^{\pi^*}(s) = V^*(s)$ for any $s \in \mathcal{S}$ where V^* is an optimal value function.

Solution. Define $f : \mathcal{S} \setminus \{(5, 1)\} \rightarrow \mathbb{N}$ as $f(s) = n$ if $n \in \mathbb{N}$ is the smallest positive integer satisfying $P^\pi(s_{t+n} = (5, 1) | s_t = s) > 0$ for some π and t . For example, $f((4, 1)) = 1, f((4, 2)) = 2, f((2, 1)) = 11$. We will show that if $f(s) = n$ then $V_m(s) = V^*(s) = \gamma^{n-1}$ for $m \geq n$. First, $V^*(s) = \gamma^{n-1}$ for $s \in f^{-1}(n)$

²The meaning of the terminal state and $R(s, a, s')$ are explained in Problem 2.

since we get a positive reward only when we reach the terminal state which implies $V^*(s) \leq \gamma^{n-1}$ and there exists an action sequence a_t, \dots, a_{t+n-1} such that $P(s_{t+n} = (5, 1) | s_t, a_t, \dots, a_{t+n-1}) = 1$ so $V^*(s) \geq \gamma^{n-1}$. (Note that the MDP is deterministic) Now, we will prove that $V_m(s) = \gamma^{f(s)-1}$ for $s \in \mathcal{S}$ such that $f(s) \leq m$ and $V_m(s) = 0$ for $s \in \mathcal{S}$ such that $f(s) > m$. For $n = 1$, $f^{-1}(1) = \{(4, 1), (5, 2)\}$. Then

$$\begin{aligned}
V_1((4, 1)) &= (LV_0)((4, 1)) = \sup_{a \in A} \mathbf{E}_{s' \sim P(\cdot | (4, 1), a)} [r((4, 1), a, s') + \gamma \cdot v(s')] \\
&= \sup_{a \in A} \mathbf{E}_{s' \sim P(\cdot | (4, 1), a)} [r((4, 1), a, s')] \\
&= \mathbf{E}_{s' \sim P(\cdot | (4, 1), E)} [r((4, 1), E, s')] = 1, \\
V_1((5, 2)) &= (LV_0)((5, 2)) = \sup_{a \in A} \mathbf{E}_{s' \sim P(\cdot | (5, 2), a)} [r((5, 2), a, s') + \gamma \cdot v(s')] \\
&= \sup_{a \in A} \mathbf{E}_{s' \sim P(\cdot | (5, 2), a)} [r((5, 2), a, s')] \\
&= \mathbf{E}_{s' \sim P(\cdot | (5, 2), S)} [r((5, 2), S, s')] = 1.
\end{aligned}$$

If $f(s) > 1$, then

$$\sup_{a \in A} \mathbf{E}_{s' \sim P(\cdot | s, a)} [r(s, a, s') + \gamma \cdot v(s')] = 0$$

since $s' \neq (5, 1)$. Suppose $V_m(s) = \gamma^{f(s)-1}$ for $s \in \mathcal{S}$ such that $f(s) \leq m$ and $V_m(s) = 0$ for $s \in \mathcal{S}$ such that $f(s) > m$. Take $s \in \mathcal{S}$. If $f(s) > m + 1$, then

$$\sup_{a \in A} \mathbf{E}_{s' \sim P(\cdot | s, a)} [r(s, a, s') + \gamma \cdot v(s')] = 0$$

since $f(s') > m$. If $f(s) \leq m + 1$, then there exists an action $a^* \in A$ such that $P(\tilde{s} | s, a^*) = 1$ for some $\tilde{s} \in \mathcal{S}$ such that $f(\tilde{s}) = f(s) - 1$ and for any $a \in A$ $f(s') \geq f(s) - 1$ where $P(s' | s, a) = 1$. Thus,

$$\sup_{a \in A} \mathbf{E}_{s' \sim P(\cdot | s, a)} [r(s, a, s') + \gamma \cdot v(s')] = r(s, a^*, \tilde{s}) + \gamma \cdot v(\tilde{s}) = \gamma^{f(s)-1}.$$

Note that $\max_{s \in \mathcal{S}} f(s) = 11$ and $\operatorname{argmax}_{s \in \mathcal{S}} f(s) = \{(2, 1)\}$. Therefore the minimum n is 10 such that all of greedy policies induced by V_n is an optimal policy since $Q(s, a)$ evaluates one more step.

4. (15pts) Let $f_{\Theta} : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ be a neural network such that

$$f_{\Theta}(\mathbf{x}) = W^{(2)} \sigma(W^{(1)} \mathbf{x} + \mathbf{b}^{(1)}) + \mathbf{b}^{(2)}$$

where

$$W^{(1)} \in \mathbb{R}^{3 \times 2}, W^{(2)} \in \mathbb{R}^{2 \times 3}, \mathbf{b}^{(1)} \in \mathbb{R}^3, \mathbf{b}^{(2)} \in \mathbb{R}^2$$

and σ is the ReLU function. Suppose the parameter $\Theta = \{W^{(1)}, W^{(2)}, \mathbf{b}^{(1)}, \mathbf{b}^{(2)}\}$ is initialized as

$$W^{(1)} = \begin{bmatrix} 1 & 0 & -1 \\ -1 & -1 & 1 \end{bmatrix}^{\top}, W^{(2)} = \begin{bmatrix} 0 & -2 & 1 \\ 1 & -1 & -1 \end{bmatrix}, \mathbf{b}^{(1)} = [0 \ 0 \ 1]^{\top}, \mathbf{b}^{(2)} = [1 \ 0]^{\top}.$$

To minimize the L^2 loss $\ell(\Theta) = \frac{1}{2} \|\mathbf{y} - f_{\Theta}(\mathbf{x})\|^2$, we will use the gradient descent method with a learning rate $\gamma = 1$. Calculate Θ for two iterations of the optimization when we have a training data

$$\mathcal{D} = \{([2 \ -3]^{\top}, [-4 \ 0]^{\top})\}.$$

To solve this problem, you cannot use the programming. Solve it by your hands.

Solution. We use the notation in the lecture slides.

$$\begin{aligned} y^0 &= \begin{bmatrix} 2 \\ -3 \end{bmatrix}, \quad x^1 = \begin{bmatrix} 1 & -1 \\ 0 & -1 \\ -1 & 1 \end{bmatrix} \begin{bmatrix} 2 \\ -3 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 5 \\ 3 \\ -4 \end{bmatrix}, \\ y^1 &= \begin{bmatrix} 5 \\ 3 \\ 0 \end{bmatrix}, \quad x^2 = \begin{bmatrix} 0 & -2 & 1 \\ 1 & -1 & -1 \end{bmatrix} \begin{bmatrix} 5 \\ 3 \\ 0 \end{bmatrix} + \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} -5 \\ 2 \end{bmatrix} \\ y^2 &= \begin{bmatrix} -5 \\ 2 \end{bmatrix} \end{aligned}$$

Also,

$$\begin{aligned} \nabla_{y^2} \ell &= \begin{bmatrix} -1 \\ 2 \end{bmatrix}, \quad \delta^2 = \begin{bmatrix} -1 \\ 2 \end{bmatrix} \odot \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} -1 \\ 2 \end{bmatrix} \\ \delta^2 (y^1)^{\top} &= \begin{bmatrix} -1 \\ 2 \end{bmatrix} [5 \ 3 \ 0] = \begin{bmatrix} -5 & -3 & 0 \\ 10 & 6 & 0 \end{bmatrix}, \quad \delta^1 = ((W^{(2)})^{\top} \delta^2) \odot \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} = \left(\begin{bmatrix} 0 & 1 \\ -2 & -1 \\ 1 & -1 \end{bmatrix} \begin{bmatrix} -1 \\ 2 \end{bmatrix} \right) \odot \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 2 \\ 0 \\ 0 \end{bmatrix}, \\ \delta^1 (y^0)^{\top} &= \begin{bmatrix} 2 \\ 0 \\ 0 \end{bmatrix} [2 \ -3] = \begin{bmatrix} 4 & -6 \\ 0 & 0 \\ 0 & 0 \end{bmatrix}. \end{aligned}$$

Thus

$$\begin{aligned} W^{(1)} &\leftarrow W^{(1)} - \delta^1 (y^0)^{\top} = \begin{bmatrix} -3 & 5 \\ 0 & -1 \\ -1 & 1 \end{bmatrix} \\ b^{(1)} &\leftarrow b^{(1)} - \delta^1 = \begin{bmatrix} -2 \\ 0 \\ 1 \end{bmatrix} \\ W^{(2)} &\leftarrow W^{(2)} - \delta^2 (y^1)^{\top} = \begin{bmatrix} 5 & 1 & 1 \\ -9 & -7 & -1 \end{bmatrix} \\ b^{(2)} &\leftarrow b^{(2)} - \delta^2 = \begin{bmatrix} 2 \\ -2 \end{bmatrix}. \end{aligned}$$

For the second iteration,

$$\begin{aligned}
y^0 &= \begin{bmatrix} 2 \\ -3 \end{bmatrix}, \quad x^1 = \begin{bmatrix} -3 & 5 \\ 0 & -1 \\ -1 & 1 \end{bmatrix} \begin{bmatrix} 2 \\ -3 \end{bmatrix} + \begin{bmatrix} -2 \\ 0 \\ 1 \end{bmatrix} = \begin{bmatrix} -23 \\ 3 \\ -4 \end{bmatrix}, \\
y^1 &= \begin{bmatrix} 0 \\ 3 \\ 0 \end{bmatrix}, \quad x^2 = \begin{bmatrix} 5 & 1 & 1 \\ -9 & -7 & -1 \end{bmatrix} \begin{bmatrix} 0 \\ 3 \\ 0 \end{bmatrix} + \begin{bmatrix} 2 \\ -2 \end{bmatrix} = \begin{bmatrix} 5 \\ -23 \end{bmatrix} \\
y^2 &= \begin{bmatrix} 5 \\ -23 \end{bmatrix}
\end{aligned}$$

Also,

$$\begin{aligned}
\nabla_{y^2} \ell &= \begin{bmatrix} 9 \\ -23 \end{bmatrix}, \quad \delta^2 = \begin{bmatrix} 9 \\ -23 \end{bmatrix} \odot \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 9 \\ -23 \end{bmatrix} \\
\delta^2(y^1)^\top &= \begin{bmatrix} 9 \\ -23 \end{bmatrix} \begin{bmatrix} 0 & 3 & 0 \end{bmatrix} = \begin{bmatrix} 0 & 27 & 0 \\ 0 & -69 & 0 \end{bmatrix}, \quad \delta^1 = ((W^2)^\top \delta^2) \odot \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} = \left(\begin{bmatrix} 5 & -9 \\ 1 & -7 \\ 1 & -1 \end{bmatrix} \begin{bmatrix} 9 & -23 \end{bmatrix} \right) \odot \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 170 \\ 0 \end{bmatrix}, \\
\delta^1(y^0)^\top &= \begin{bmatrix} 0 \\ 170 \\ 0 \end{bmatrix} \begin{bmatrix} 2 & -3 \end{bmatrix} = \begin{bmatrix} 0 & 0 \\ 340 & -510 \\ 0 & 0 \end{bmatrix}.
\end{aligned}$$

Thus

$$\begin{aligned}
W^{(1)} &\leftarrow W^{(1)} - \delta^1(y^0)^\top = \begin{bmatrix} -3 & 5 \\ -340 & 509 \\ -1 & 1 \end{bmatrix} \\
b^{(1)} &\leftarrow b^{(1)} - \delta^1 = \begin{bmatrix} -2 \\ -170 \\ 1 \end{bmatrix} \\
W^{(2)} &\leftarrow W^{(2)} - \delta^2(y^1)^\top = \begin{bmatrix} 5 & -26 & 1 \\ -9 & 62 & -1 \end{bmatrix} \\
b^{(2)} &\leftarrow b^{(2)} - \delta^2 = \begin{bmatrix} -7 \\ 21 \end{bmatrix}.
\end{aligned}$$