

$$Q3-(c) \quad V_{\pi} = \sum_a \pi(a|s) \sum_{s',r} P(s',r|s,a) [r + \gamma V_{\pi}(s')]$$

For A: $s' \in \{B, t\}$ ($t \rightarrow \text{terminal}$) $V_{\pi}(A) = \pi(\rightarrow, A) \sum_{s',r} P(s',r|A,\rightarrow) [r + \gamma V_{\pi}(s')] + \pi(\leftarrow, A) \sum_{s',r} P(s',r|A,\leftarrow) [r + \gamma V_{\pi}(s')]$

$$\begin{aligned} \rightarrow V_{\pi_A} &= 0.9 \times [P(B,r|A,\rightarrow)(r + \gamma V_{\pi}(B)) + P(t,r|A,\rightarrow)(r + \gamma V_{\pi}(t))] \\ &\quad + 0.1 [P(B,r|A,\leftarrow)(r + \gamma V_{\pi}(B)) + P(t,r|A,\leftarrow)(r + \gamma V_{\pi}(t))] = 0.9 \times 1 \times (0 + \gamma V_{\pi}(B)) \\ &\quad + 0.1 \times 1 \times (100) = 0.9 V_{\pi}(B) + 10 \quad (1) \end{aligned}$$

For B: $s' \in \{A, C\}$ $V_{\pi}(B) = \pi(\rightarrow, B) \sum_{s',r} P(s',r|B,\rightarrow) [r + \gamma V_{\pi}(s')] + \pi(\leftarrow, B) \sum_{s',r} P(s',r|B,\leftarrow) [r + \gamma V_{\pi}(s')]$

$$\begin{aligned} \rightarrow V_{\pi_B} &= 0.9 [P(C,r|B,\rightarrow)(r + \gamma V_{\pi}(C)) + P(A,r|B,\rightarrow)(r + \gamma V_{\pi}(A))] \\ &\quad + 0.1 [P(C,r|B,\leftarrow)(r + \gamma V_{\pi}(C)) + P(A,r|B,\leftarrow)(r + \gamma V_{\pi}(A))] = 0.9 \times 1 \times (0 + \gamma V_{\pi}(C)) \\ &\quad + 0.1 \times 1 \times (0 + \gamma V_{\pi}(A)) = 0.9 V_{\pi}(C) + 0.1 V_{\pi}(A) \quad (2) \end{aligned}$$

For C: $s' \in \{B, t\} \Rightarrow V_{\pi}(C) = \pi(\rightarrow, C) \sum_{s',r} P(s',r|C,\rightarrow) [r + \gamma V_{\pi}(s')] + \pi(\leftarrow, C) \sum_{s',r} P(s',r|C,\leftarrow) [r + \gamma V_{\pi}(s')]$

$$\begin{aligned} \rightarrow V_{\pi_C} &= 0.9 [P(t,r|C,\rightarrow)(r + \gamma V_{\pi}(t)) + P(B,r|C,\rightarrow)(r + \gamma V_{\pi}(B))] \\ &\quad + 0.1 [P(t,r|C,\leftarrow)(r + \gamma V_{\pi}(t)) + P(B,r|C,\leftarrow)(r + \gamma V_{\pi}(B))] = 0.9 \times 1 \times (-1 + \gamma \times 0) + 0.1 \times 1 \times (0 + \gamma V_{\pi}(B)) \\ &\quad = -0.9 + 0.1 V_{\pi}(B) \quad (3) \end{aligned}$$

①, ②, ③ $V_{\pi}(A) \approx 10.2085$ $V_{\pi}(B) \approx 0.2317$ $V_{\pi}(C) \approx -0.8768$

(b) For A: $\max Q(s,a) = 2$ (for left) $\xrightarrow{\epsilon\text{-greedy}}$ left: $\frac{1}{4} \times 0.25 + 0.75 = 0.8125$
up: $\frac{1}{4} \times 0.25 = 0.0625$
down: $\frac{1}{4} \times 0.25 = 0.0625$
right: $\frac{1}{4} \times 0.25 = 0.0625$

For B: $\max Q(s,A) = 1.2$ (for right)

$\xrightarrow{\epsilon\text{-greedy}}$ left: $\frac{1}{4} \times 0.25 = 0.0625$
up: $\frac{1}{4} \times 0.25 = 0.0625$
down: $\frac{1}{4} \times 0.25 = 0.0625$
right: $\frac{1}{4} \times 0.25 + 0.75 = 0.8125$

$$\Rightarrow$$

state	up	down	left	right
A	0.0625	0.0625	0.8125	0.0625
B	0.0625	0.0625	0.0625	0.8125

$$d) \quad Q(A, up) \leftarrow Q(A, up) + \alpha [R_{t+1} + \gamma (\sum_a \pi(a|B) Q(B, a)) - Q(A, up)]$$

we have: $Q(A, up) = -1.2$, $R_{t+1} = -1$, $\gamma = 0.99$, $\alpha = 0.1$ $a \in \{up, down, left, right\}$

$$\rightarrow Q(A, up) \leftarrow -1.2 + 0.1 [-1 + 0.99 (0.0625 \times (-0.2 + 0.1 - 1.1) + 0.8125 \times (1.2)) + 1.2] = -1.0909$$

Q5: (a) A particular part of a city might have had a bad reputation regarding crime, because of poverty or a culture. But even if these circumstances change over the years, and the next generation become good people, the historical data of their home location predicts a higher chance of being criminal for those people. It somehow even relates to the bias towards black people.

(b) The algorithm relies on input data from human. Some feature that the algorithm learns, are biased, since the dataset was imbalanced. For example, the network might not have observed enough input data from minorities. Or in some cases, the algorithm itself might not do well on some races and skin colors.