

RL Tutorial worksheet

Exercise 1 – Q iteration:

$$Q_{new}(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha(r_t + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t))$$

$$Q_0 = \begin{matrix} & \begin{matrix} l & w \end{matrix} \\ \begin{matrix} H \\ N \\ F \\ Q \end{matrix} & \begin{bmatrix} 10 & 8 \\ 10 & 8 \\ 10 & 8 \\ 0 & 0 \end{bmatrix} \end{matrix} \quad \begin{matrix} \alpha = 0.8 \\ \gamma = 0.9 \end{matrix}$$

User 1: $N \xrightarrow{l} F \xrightarrow{l} Q$:

<p>Observation 1:</p> <p>$(s_t, a_t, r_t, s_{t+1}) = (N, l, +1, F)$</p> <p>$Q_{new}(\dots, \dots) \leftarrow Q(\dots, \dots) + \alpha(\dots + \gamma \max_a Q(\dots, a) - Q(\dots, \dots))$</p>	$Q_1 = \begin{matrix} & \begin{matrix} l & w \end{matrix} \\ \begin{matrix} H \\ N \\ F \\ Q \end{matrix} & \begin{bmatrix} \square & \square \\ \square & \square \\ \square & \square \\ 0 & 0 \end{bmatrix} \end{matrix}$
<p>Observation 2:</p> <p>$(s_t, a_t, r_t, s_{t+1}) = (\dots, \dots, \dots, \dots)$</p> <p>$Q_{new}(\dots, \dots) \leftarrow Q(\dots, \dots) + \alpha(\dots + \gamma \max_a Q(\dots, a) - Q(\dots, \dots))$</p>	$Q_2 = \begin{matrix} & \begin{matrix} l & w \end{matrix} \\ \begin{matrix} H \\ N \\ F \\ Q \end{matrix} & \begin{bmatrix} \square & \square \\ \square & \square \\ \square & \square \\ 0 & 0 \end{bmatrix} \end{matrix}$

User 2: $H \xrightarrow{l} N \xrightarrow{l} N \xrightarrow{l} F \xrightarrow{w} N$:

<p>Observation 3:</p> <p>$(s_t, a_t, r_t, s_{t+1}) = (\dots, \dots, \dots, \dots)$</p> <p>$Q_{new}(\dots, \dots) \leftarrow Q(\dots, \dots) + \alpha(\dots + \gamma \max_a Q(\dots, a) - Q(\dots, \dots))$</p>	$Q_3 = \begin{array}{ c c c } \hline & l & w \\ \hline H & \square & \square \\ \hline N & \square & \square \\ \hline F & \square & \square \\ \hline Q & 0 & 0 \\ \hline \end{array}$
<p>Observation 4:</p> <p>$(s_t, a_t, r_t, s_{t+1}) = (\dots, \dots, \dots, \dots)$</p> <p>$Q_{new}(\dots, \dots) \leftarrow Q(\dots, \dots) + \alpha(\dots + \gamma \max_a Q(\dots, a) - Q(\dots, \dots))$</p>	$Q_4 = \begin{array}{ c c c } \hline & l & w \\ \hline H & \square & \square \\ \hline N & \square & \square \\ \hline F & \square & \square \\ \hline Q & 0 & 0 \\ \hline \end{array}$
<p>Observation 5:</p> <p>$(s_t, a_t, r_t, s_{t+1}) = (\dots, \dots, \dots, \dots)$</p> <p>$Q_{new}(\dots, \dots) \leftarrow Q(\dots, \dots) + \alpha(\dots + \gamma \max_a Q(\dots, a) - Q(\dots, \dots))$</p>	$Q_5 = \begin{array}{ c c c } \hline & l & w \\ \hline H & \square & \square \\ \hline N & \square & \square \\ \hline F & \square & \square \\ \hline Q & 0 & 0 \\ \hline \end{array}$
<p>Observation 6:</p> <p>$(s_t, a_t, r_t, s_{t+1}) = (\dots, \dots, \dots, \dots)$</p> <p>$Q_{new}(\dots, \dots) \leftarrow Q(\dots, \dots) + \alpha(\dots + \gamma \max_a Q(\dots, a) - Q(\dots, \dots))$</p>	$Q_6 = \begin{array}{ c c c } \hline & l & w \\ \hline H & \square & \square \\ \hline N & \square & \square \\ \hline F & \square & \square \\ \hline Q & 0 & 0 \\ \hline \end{array}$

Exercise 2 – Value Iteration:

$$V_{new}(s) \leftarrow \max_a \left\{ \sum_{s'} P(s'|s, a) (r(s, a, s') + \gamma V(s')) \right\}, \quad \gamma = 0.9$$

$$V_0 = \begin{matrix} H \\ N \\ F \\ Q \end{matrix} \begin{bmatrix} 10 \\ 10 \\ 10 \\ 0 \end{bmatrix}$$

	Next State				
		H	N	F	Q
Current State	H	0.8	0.1	0.0	0.1
	N	0.0	0.5	0.2	0.3
	F	0.0	0.0	0.2	0.8
When user loses: a = l					

	Next State				
		H	N	F	Q
Current State	H	0.9	0.0	0.0	0.1
	N	0.7	0.2	0.0	0.1
	F	0.3	0.4	0.1	0.2
When user wins: a = w					

$$V_{new}(F)$$

$$\leftarrow \max \left\{ \begin{aligned} & \dots \times (\dots + \gamma \times \dots) + \dots \times (\dots + \gamma \times \dots) + \dots \times (\dots + \gamma \times \dots) + P(Q|F, l) \times (\dots + \gamma \times \dots) \\ & \dots \times (\dots + \gamma \times \dots) + \dots \times (\dots + \gamma \times \dots) + \dots \times (\dots + \gamma \times \dots) + P(Q|F, w) \times (\dots + \gamma \times \dots) \end{aligned} \right\}$$

$$= \dots$$

$$V_{new}(N)$$

$$\leftarrow \max \left\{ \begin{aligned} & P(H, N, l) \times (\dots + \gamma \times \dots) + \dots \times (\dots + \gamma \times \dots) + \dots \times (\dots + \gamma \times \dots) + \dots \times (\dots + \gamma \times \dots) \\ & P(H, N, w) \times (\dots + \gamma \times \dots) + \dots \times (\dots + \gamma \times \dots) + \dots \times (\dots + \gamma \times \dots) + \dots \times (\dots + \gamma \times \dots) \end{aligned} \right\}$$

$$= \dots$$

$$V_{new}(H)$$

$$\leftarrow \max \left\{ \begin{aligned} & \dots \times (\dots + \gamma \times \dots) + \dots \times (\dots + \gamma \times \dots) + \dots \times (\dots + \gamma \times \dots) + \dots \times (\dots + \gamma \times \dots) \\ & \dots \times (\dots + \gamma \times \dots) + \dots \times (\dots + \gamma \times \dots) + \dots \times (\dots + \gamma \times \dots) + \dots \times (\dots + \gamma \times \dots) \end{aligned} \right\}$$

$$= \dots$$