

Q-Learning

Shusen Wang

Sarsa VS Q-Learning

- Sarsa is for training action-value function, $Q_{\pi}(s, a)$.
- TD target: $y_t = r_t + \gamma \cdot Q_{\pi}(s_{t+1}, a_{t+1})$.
- We used Sarsa for updating value network (critic).

Sarsa VS Q-Learning

- Q-learning is for training the optimal action-value function, $Q^*(s, a)$.
- TD target: $y_t = r_t + \gamma \cdot \max_a Q^*(s_{t+1}, a)$.
- We used Q-learning for updating DQN.

Derive TD Target

Derive TD Target

- We have proved that for all π ,

$$Q_{\pi}(s_t, a_t) = \mathbb{E}[R_t + \gamma \cdot Q_{\pi}(S_{t+1}, A_{t+1})].$$

Derive TD Target

- We have proved that for all π ,

$$Q_{\pi}(s_t, a_t) = \mathbb{E}[R_t + \gamma \cdot Q_{\pi}(S_{t+1}, A_{t+1})].$$

- If π is the optimal policy π^* , then

$$Q_{\pi^*}(s_t, a_t) = \mathbb{E}[R_t + \gamma \cdot Q_{\pi^*}(S_{t+1}, A_{t+1})].$$

Derive TD Target

- We have proved that for all π ,

$$Q_{\pi}(s_t, a_t) = \mathbb{E}[R_t + \gamma \cdot Q_{\pi}(S_{t+1}, A_{t+1})].$$

- If π is the optimal policy π^* , then

$$Q_{\pi^*}(s_t, a_t) = \mathbb{E}[R_t + \gamma \cdot Q_{\pi^*}(S_{t+1}, A_{t+1})].$$

- Q_{π^*} and Q^* both denote *the optimal action-value function*.

Derive TD Target

- We have proved that for all π ,

$$Q_{\pi}(s_t, a_t) = \mathbb{E}[R_t + \gamma \cdot Q_{\pi}(S_{t+1}, A_{t+1})].$$

- If π is the optimal policy π^* , then

$$Q_{\pi^*}(s_t, a_t) = \mathbb{E}[R_t + \gamma \cdot Q_{\pi^*}(S_{t+1}, A_{t+1})].$$

- Q_{π^*} and Q^* both denote *the optimal action-value function*.

Identity: $\underline{Q^*(s_t, a_t)} = \mathbb{E}[\underline{R_t + \gamma \cdot Q^*(S_{t+1}, A_{t+1})}]$.

Derive TD Target

Identity: $Q^*(s_t, a_t) = \mathbb{E}[R_t + \gamma \cdot Q^*(s_{t+1}, A_{t+1})]$.

- The action A_{t+1} is computed by

$$A_{t+1} = \underset{a}{\operatorname{argmax}} Q^*(s_{t+1}, a).$$

- Thus $Q^*(s_{t+1}, \underline{A_{t+1}}) = \max_a Q^*(s_{t+1}, a)$.

Derive TD Target

Identity: $Q^*(s_t, a_t) = \mathbb{E}[R_t + \gamma \cdot Q^*(S_{t+1}, A_{t+1})]$.

- Thus $Q^*(S_{t+1}, A_{t+1}) = \max_a Q^*(S_{t+1}, a)$.

Derive TD Target

Identity: $Q^*(s_t, a_t) = \mathbb{E}[R_t + \gamma \cdot Q^*(S_{t+1}, A_{t+1})]$.

$$= \max_a Q^*(S_{t+1}, a)$$




Identity: $Q^*(s_t, a_t) = \mathbb{E}\left[R_t + \gamma \cdot \max_a Q^*(S_{t+1}, a)\right]$.

Derive TD Target

Identity: $Q^*(s_t, a_t) = \mathbb{E} \left[R_t + \gamma \cdot \max_a Q^*(S_{t+1}, a) \right].$

Derive TD Target

Identity: $Q^*(s_t, a_t) = \mathbb{E} \left[R_t + \gamma \cdot \max_a Q^*(s_{t+1}, a) \right].$

- Let (s_{t+1}, r_t) be an observation of (S_{t+1}, R_t) .
- TD target: $y_t = r_t + \gamma \cdot \max_a Q^*(s_{t+1}, a).$


Derive TD Target

Identity: $Q^*(s_t, a_t) = \mathbb{E} \left[R_t + \gamma \cdot \max_a Q^*(s_{t+1}, a) \right]$

$\approx y_t$

- Let (s_{t+1}, r_t) be an observation of (S_{t+1}, R_t) .
- TD target: $y_t = r_t + \gamma \cdot \max_a Q^*(s_{t+1}, a)$.

Q-Learning: Tabular Version

Q-Learning (tabular version)

- Observe (s_t, a_t, r_t, s_{t+1}) .
- TD target: $y_t = r_t + \gamma \cdot \max_a Q^*(s_{t+1}, a)$.

Q-Learning (tabular version)

- Observe (s_t, a_t, r_t, s_{t+1}) .
- TD target: $y_t = r_t + \gamma \max_a Q^*(s_{t+1}, a)$.

	Action a_1	Action a_2	Action a_3	Action a_4	...
State s_1					
State s_2					
State s_3					
\vdots					

Q-Learning (tabular version)

- Observe (s_t, a_t, r_t, s_{t+1}) .
- TD target: $y_t = r_t + \gamma \cdot \max_a Q^*(s_{t+1}, a)$.
- TD error: $\delta_t = Q^*(s_t, a_t) - y_t$.
- Update: $Q^*(s_t, a_t) \leftarrow Q^*(s_t, a_t) - \alpha \cdot \delta_t$.

make $Q^*(s_t, a_t)$ closer to y_t

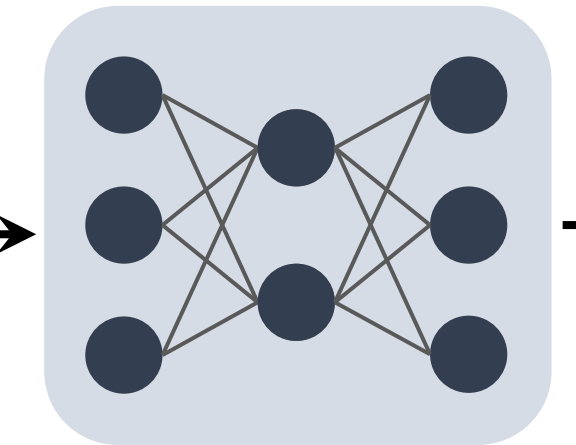
Q-Learning: DQN Version

DQN Version

- Approximate $Q^*(s, a)$ by DQN, $Q(s, a|w)$.



state s



DQN
(parameterized by w)



$Q(s, \text{"left"}; w)$

$Q(s, \text{"right"}; w)$

$Q(s, \text{"up"}; w)$

DQN Version

- Approximate $Q^*(s, a)$ by DQN, $Q(s, a|\mathbf{w})$.
- DQN controls the agent by: $a_t = \underset{a}{\operatorname{argmax}} Q(s_t, a|\mathbf{w})$
- We seek to learn the parameter, \mathbf{w} .

Q-Learning (DQN Version)

- Observe (s_t, a_t, r_t, s_{t+1}) .
- TD target: $y_t = r_t + \gamma \cdot \max_a Q(s_{t+1}, a | \mathbf{w})$.
- TD error: $\delta_t = Q(s_t, a_t | \mathbf{w}) - y_t$.
- Update: $\mathbf{w} \leftarrow \mathbf{w} - \alpha \cdot \delta_t \cdot \frac{\partial q(s_t, a_t | \mathbf{w})}{\partial \mathbf{w}}$.

Summary

- **Goal:** Learn the optimal action-value function Q^* .
- **Tabular version** (directly learn Q^*).
 - There are finite states and actions.
 - Draw a table, and update the table by Q-learning.
- **DQN version** (function approximation).
 - Approximate Q^* by the DQN, $Q(\textcolor{green}{s}, \textcolor{red}{a}|\mathbf{w})$.
 - Update the parameter, \mathbf{w} , by Q-learning.

Thank you!