

DQN

Goal: Approximate the optimal action-value function

$$Q^*(s, a) = \max_{\pi} \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_t \mid s_0 = s, a_0 = a, \pi \right]$$

$$Q^*(s, a) = \mathbb{E} [r + \gamma \max_{a'} Q(s', a') \mid s, a] - \text{Bellman optimality condition}$$

Classical Q-Learning and the Bellman equation

$$Q(s, a) \leftarrow Q(s, a) + \alpha [r + \gamma \max_{a'} Q(s', a') - Q(s, a)]$$

↑ Bellman equation

$$\text{If } \alpha = 1 \rightarrow Q(s, a) \leftarrow r + \gamma \max_{a'} Q(s', a')$$

Approximation with a NN

$$Q(s, a; \theta) \approx Q^*(s, a)$$

Stabilize learning: Two-Network Architecture

Challenges: Instability due to non-stationary targets

Policy Network: Estimate Q-values for current state

- Updated continuously at every training step

Target Network: Provides stable target Q-values during the training of the Policy Network.

- Copy of Policy Network but updated less frequently

Implementation:

- Action selection: Policy Network (state)
- Environment interaction: Action executed $\mapsto s', r$
- Experience storage: (s, a, r, s') stored in replay buffer \uparrow Experience replay
- Sampling & Training:
 - * - Sample from replay buffer - Avoids overfitting, improves generalization
 - Compute target Q-value: $y = r + \gamma \max_{a'} Q_{\text{target}}(s', a')$
- Updating the Policy Network \rightarrow Update with gradient
 - $L = [y - Q_{\text{policy}}(s, a)]^2$ (Loss) dextra (Lecture 12.5)

* $q\text{-value} = \text{policy_net}(\text{shape})$ - dimension (batch-size, actions)
 $B_1 [[24, 12]]$
 $B_2 [[13, 1]]$
 $\text{policy_net.gather}(1, \text{actions})$
 $\Rightarrow [[24], [13]]$
 $\text{actions} = [[0], [0]]$
 $\text{squeeze}() \Rightarrow [24, 13] = q\text{-value}$

target-net (next-state)

$$b_1 \begin{bmatrix} a_1 & a_2 \\ 2 & 0 \end{bmatrix}, \\ b_1 \begin{bmatrix} 0 & 3 \end{bmatrix}$$

target-net.max(1) \mapsto ([2,3], [0,1])

[0] \mapsto [2,3] = ~~target~~ ^{next} q-value if done = True = 1

$$\text{target-q-value} = r + \gamma \cdot \text{next-q-value} \cdot (1 - \text{done})$$

Recall

$$\underbrace{Q(s,a)}_{\text{Policy-net}} = r + \gamma \max_{a'} \underbrace{Q(s',a')}_{\text{Target-net}} \quad (\text{Bellman equation})$$

$$\Rightarrow \text{Loss} = |\text{target-q-value} - \text{q-value}|^2$$