

# Zadanie 6

## Q-Learning

Kacper Kania

28.04.2024

### Treść zadania

Q-Learning jest podstawowym algorytmem uczenia ze wzmocnieniem. Jest to metoda uczenia bez nadzoru, która polega na uczeniu się strategii zachowań w środowisku, aby osiągnąć maksymalną nagrodę. W algorytmie Q-Learningu agent uczy się funkcji wartości akcji, która określa wartość każdej akcji w każdym stanie. W każdym kroku agent wybiera akcję, która maksymalizuje wartość funkcji wartości akcji. Wartość funkcji wartości akcji jest aktualizowana na podstawie nagrody otrzymanej po wykonaniu akcji i wartości funkcji wartości akcji w nowym stanie. Państwo zrealizujecie funkcję Q jako tablicę, natomiast Q może być dowolną funkcją, nawet wyuczalną. Na przykład, jedna ze słynniejszych prac DeepMindu wykorzystywała głębokie sieci neuronowe do nauczenia funkcji Q z pikseli obrazu gry [1].

Państwa zadaniem jest implementacja uczenia funkcji Q jako tablicy w środowisku **FrozenLake-v1**<sup>1</sup> z biblioteki **gym** z ustawieniami:

```
gym.make('FrozenLake-v1', desc=None, map_name="8x8", is_slippery=True)
```

Sama biblioteka zwraca dla podanej akcji obserwację (zmianę stanu środowiska), nagrodę oraz informację o zakończeniu epizodu.

```
import gym
env = gym.make("LunarLander-v2", render_mode="human")
observation, info = env.reset(seed=42)
for _ in range(1000):
    action = policy(observation) # User-defined policy function
    observation, reward, terminated, truncated, info = env.step(action)

    if terminated or truncated:
        observation, info = env.reset()
env.close()
```

Dla przypomnienia jeden krok uczenia Q wygląda następująco:

$$Q^{\text{new}}(S_t, A_t) \leftarrow \left( 1 - \underbrace{\alpha}_{\text{learning rate}} \right) \cdot Q(S_t, A_t) + \alpha \cdot \left( \underbrace{R_{t+1}}_{\text{reward}} + \underbrace{\gamma}_{\text{discount factor}} \cdot \underbrace{Q(S_{t+1}, a)}_{\text{Q value in next state after selecting action according to the policy}} \right), \quad (1)$$

new value

gdzie:

<sup>1</sup>[https://www.gymnasium.dev/environments/toy\\_text/frozen\\_lake/](https://www.gymnasium.dev/environments/toy_text/frozen_lake/)

- $\alpha$  to współczynnik uczenia,
- $R_{t+1}$  to nagroda otrzymana po wykonaniu akcji  $A_t$  w stanie  $S_t$ ,
- $\gamma$  to współczynnik dyskontowania,
- $Q(\cdot, \cdot)$  nasza funkcja do wyuczenia,
- $\pi_{a(\cdot)}$  zwraca akcję  $a$  w zależności od określonej polityki wyboru i wartości  $Q$  (np. polityka  $\epsilon$ -zachłanna).

W sprawozdaniu należy:

- zbadać wpływ parametru uczenia  $\alpha$  na zbieżność algorytmu,
- zbadać wpływ parametru dyskontowania  $\gamma$  na zbieżność algorytmu,
- zbadać wpływ parametru eksploracji  $\epsilon$  na zbieżność algorytmu w podejściu  $\epsilon$ -zachłannym, oraz parametr  $T$  w strategii opartej na rozkładzie Boltzmanna.

Zbieżność algorytmu proszę przedstawić jako wykres zależności końcowej nagrody od kroku uczenia funkcji  $Q$ . W tabelach natomiast proszę przestawić średnią nagrodę w ostatnich 10 krokach uczenia wraz z odchyleniem standardowym oraz liczbę rozgrywek zakończonych sukcesem (10 rozgrywek). Proszę przedstawić wyniki dla 5 różnych random seedów.

## Bibliografia

- [1] V. Mnih *et al.*, “Playing atari with deep reinforcement learning,” *arXiv preprint arXiv:1312.5602*, 2013.