

Raport

Alesia Filinkova

336180

1 Treść ćwiczenia

Zaimplementować algorytm Q-learning, a następnie użyć go do wytrenowania agenta rozwiązującego problem Cliff Walking

(https://gymnasium.farama.org/environments/toy_text/cliff_walking/).

Stworzyć wizualizacje wyuczonej polityki

2 Doprecyzowanie

- Środowisko CliffWalking-v0 zawiera plansze o wymiarach 4x12, gdzie celem jest dotarcie z pola startowego (w lewym dolnym rogu) do celu (w prawym dolnym rogu) bez upadku z klifu.
- Użyto tablicy Q table, aby przechowywać wartość oczekiwanych nagród dla każdej kombinacji stanu i akcji.
- Metryka oceny jest suma nagród uzyskanych podczas treningu oraz wizualizacja wyuczonej polityki.

3 Cel i opis eksperymentów

3.1 Cel

Celem ćwiczenia było:

1. Zaimplementowanie algorytmu Q-learning.
2. Wytrenowanie agenta, który skutecznie nauczy się omijać klif i osiągnie stan końcowy z minimalnym kosztem.
3. Wizualizacja wyników i analiza zachowania agenta.

3.2 Opis eksperymentów:

- Zbiór danych: Plansza CliffWalking (4x12), dostępna w gymnasium.
- Algorytm: Q-learning z tablicą q
- Metryka: Wizualizacja polityki, dystrybucje odwiedzanych stanów i wykonanych akcji
- Parametry:
 - Współczynnik uczenia (learning rate)
 - Współczynnik dyskontowania (discount factor)
 - Prawdopodobieństwo eksploracji (epsilon)
 - Liczba epizodów treningowych (num episodes)

4 Przygotowanie środowiska i danych

Skrypt można uruchomić przez terminal za pomocą polecenia:

1. `git clone https://gitlab-stud.elka.pw.edu.pl/afilinko/wsi.git`
2. `python3 -m venv venv`
3. `source venv/bin/activate`
4. `cd /lab6`
5. `pip install -r requirements.txt`
6. `python3 main.py`

5 Wyniki

Eksperymentalnie ustawiono, że najlepsze wyniki dają następujące wartości hiperparametrów: $learning_rate = 0,15$, $gamma = 0,95$, $epsilon = 0,1$, $episodes = 500$.

Można ustawić większą liczbę iteracji, jednak to nie bardzo polepsza wynik, skoro cały algorytm jest oparty o dość losowe zdarzenia.

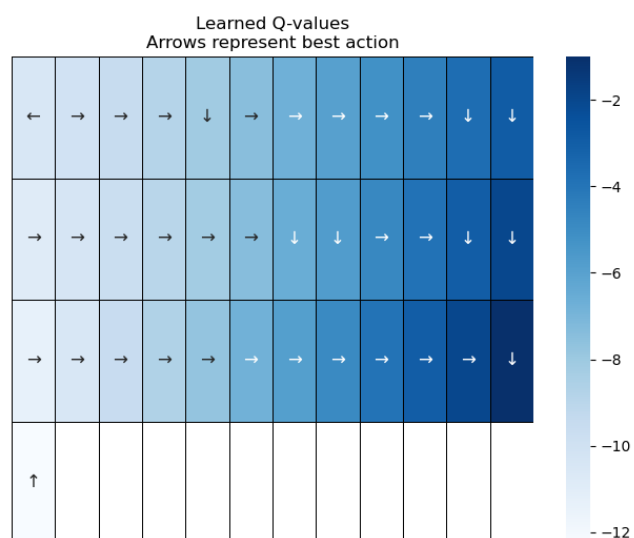


Figure 1: Wizualizacja wyuczonej polityki

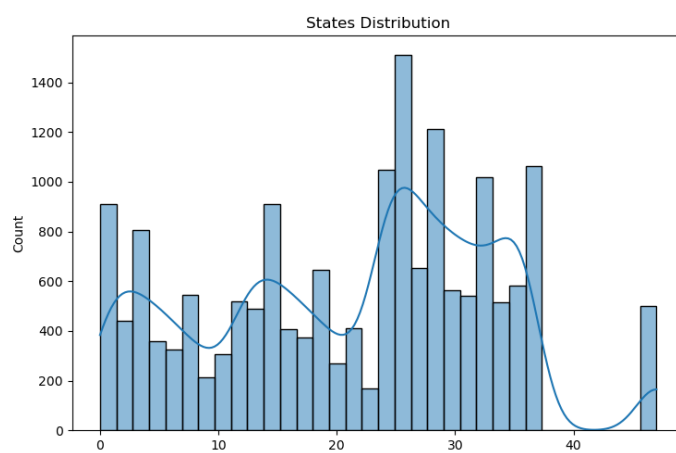
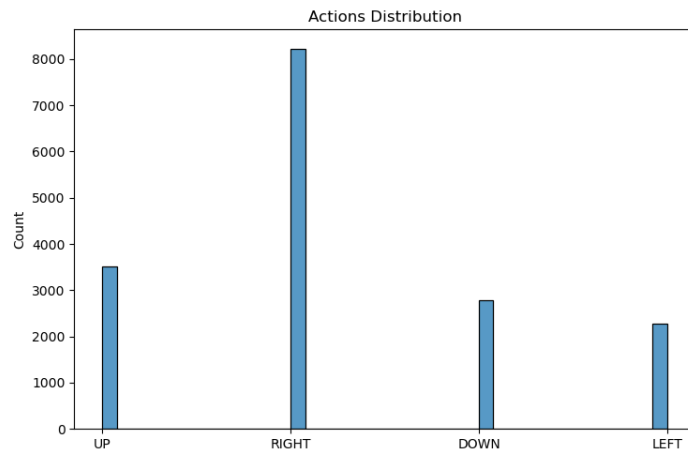


Figure 2: Ilość kroków



6 Wnioski

- Agent skutecznie nauczył się omijać klif, poruszając się wzdłuż bezpiecznych pól, co widać po spójnej polityce (\rightarrow w pierwszym rzędzie, a następnie \downarrow w celu dotarcia do końca).
- Liczba odwiedzanych stanów i akcji wskazuje, że agent w większości epizodów trzymał się bezpiecznych pól
- Wizualizacja polityki pokazuje, że agent preferuje ruchy w kierunku celu, a rozkłady stanów i akcji wskazują na efektywne eksplorowanie środowiska.
- W przyszłości można by zastosować inne techniki uczenia, takie jak Deep Q-Learning, aby zbadać ich wpływ na wyniki.