

WSTĘP DO SZTUCZNEJ INTELIGENCJI

Ćwiczenie 5 – Algorytm Q-learning

JAKUB KWAŚNIAK 331396

Wstęp

1. Zaimplementować algorytm Q-learning, a następnie użyć go do wytrenowania agenta rozwiązującego problem Cliff Walking https://gymnasium.farama.org/environments/toy_text/cliff_walking/
2. Stworzyć wizualizację wyuczonej polityki i umieścić ją w sprawozdaniu. Wzór wizualizacji https://gymnasium.farama.org/tutorials/training_agents/FrozenLake_tuto/#visualization

1. Zaimplementowany został zachłanny algorytm Q-learning (Epsilon-Greedy), który wyznacza najlepszą dla siebie akcję w danym stanie według formuły Bellmana:

Bellman's Equation:

Update $Q(s,a) := Q(s,a) + lr [R(s,a) + \gamma \max_{a'} Q(s',a') - Q(s,a)]$

$\Delta = [R(s,a) + \gamma \max_{a'} Q(s',a') - Q(s,a)]$

$Q(s,a) := Q(s,a) + lr * \Delta$

w prostych przypadkach takich jak cliff walking często w pierwszych epizodach (kiedy świat nie jest jeszcze oceniony) $\max_{a'} Q(s',a') == Q(s,a)$ więc gamma nie powinna być równa 1

W zaimplementowanym środowisku z pakietu gymnasium nagrody są następujące:

Każdy krok -> -1

Dojście do celu -> 0

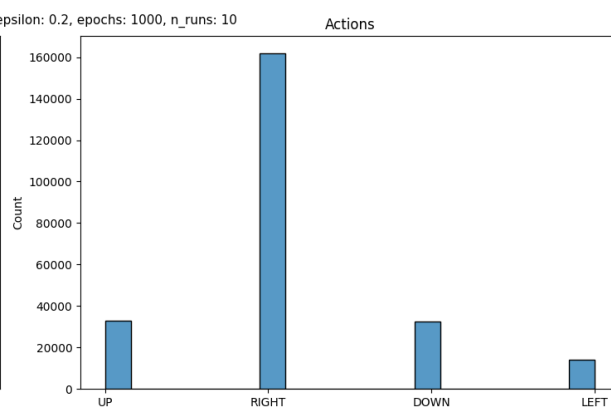
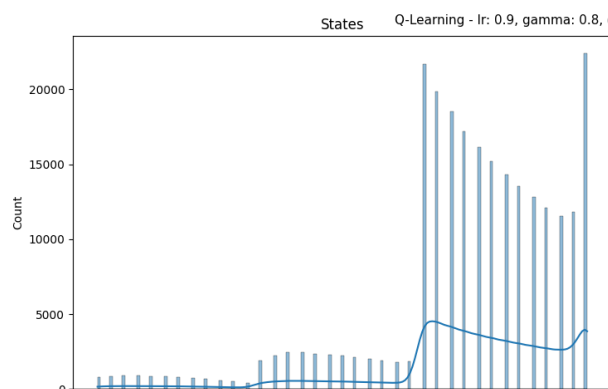
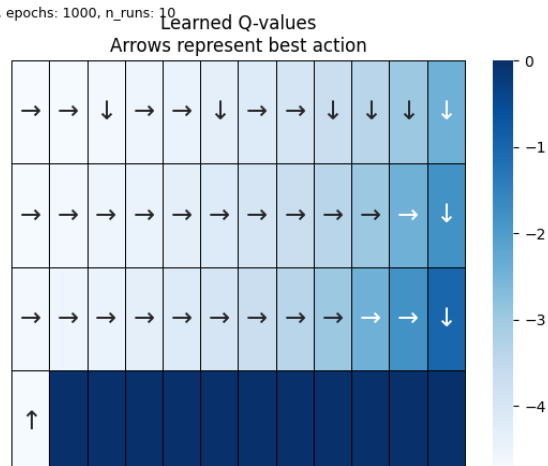
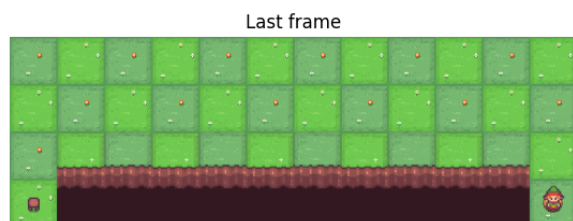
Spadnięcie z klifu -> -100

Dlatego ustawienie $\gamma=1$ uniemożliwiłoby algorytmowi ocenę środowiska i naukę dobrej taktyki

Użycie Epsilon – Greedy pozwala w pewnym stopniu rozwiązać problem „Exploration-exploitation dilemma” – jako parametr epsilon przyjmuje wartość będącą prawdopodobieństwem eksploracji, zwyczajowo jest to mała wartość by uczeń mógł eksploatować najlepszą znaną część środowiska a jednocześnie czasami eksplorować resztę środowiska i nie zamykał się na inne możliwości – może znaleźć równie dobrą lub lepszą część środowiska.

Eksperyment 1:

Q-Learning - lr: 0.9, gamma: 0.8, epsilon: 0.2, epochs: 1000, n runs: 10



Qtable

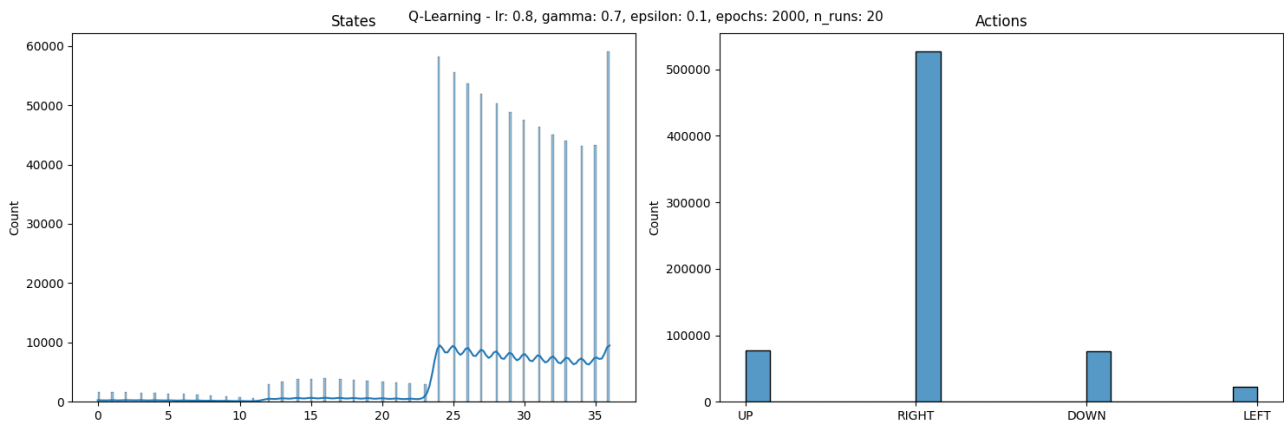
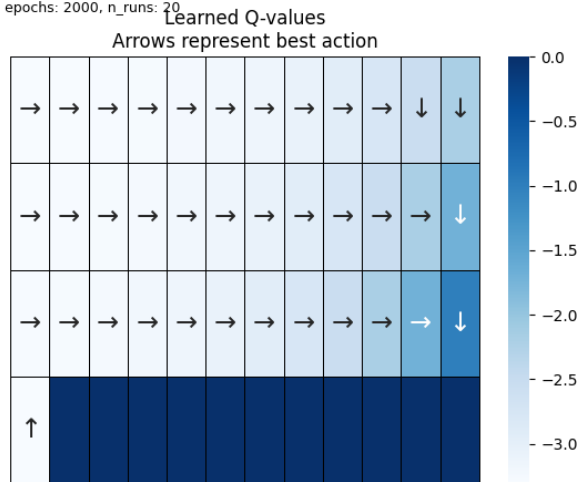
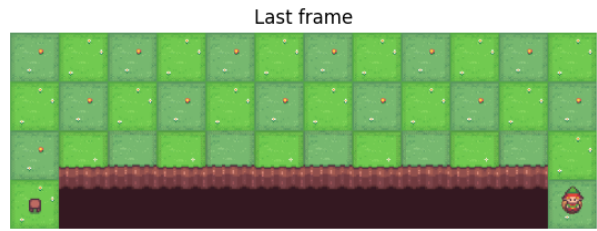
Cumulative Rewards

[illegible]

	Episodes	Rewards	Steps	cum_rewards	map_size
0	0	-866.0	173.0	-866.0	4x12
1	1	-1437.0	447.0	-2303.0	4x12
2	2	-653.0	158.0	-2956.0	4x12
3	3	-275.0	77.0	-3231.0	4x12
4	4	-44.0	44.0	-3275.0	4x12
...
9995	995	-16.0	16.0	-108151.0	4x12
9996	996	-19.0	19.0	-108170.0	4x12
9997	997	-13.0	13.0	-108183.0	4x12
9998	998	-13.0	13.0	-108196.0	4x12
9999	999	-238.0	40.0	-108434.0	4x12

Eksperyment 2:

Q-Learning - lr: 0.8, gamma: 0.7, epsilon: 0.1, epochs: 2000, n_runs: 20



Qtable

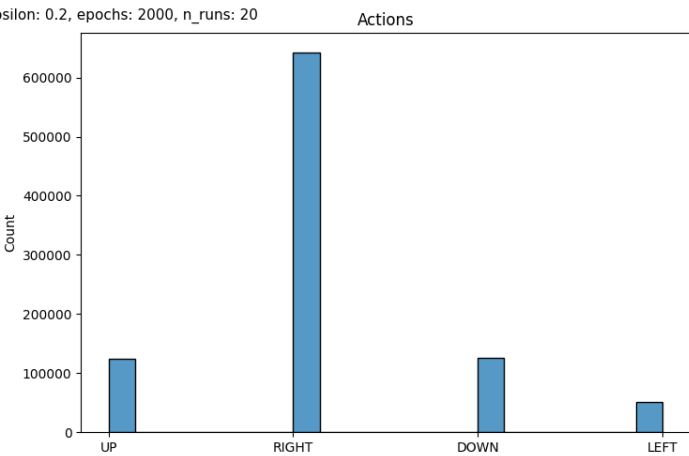
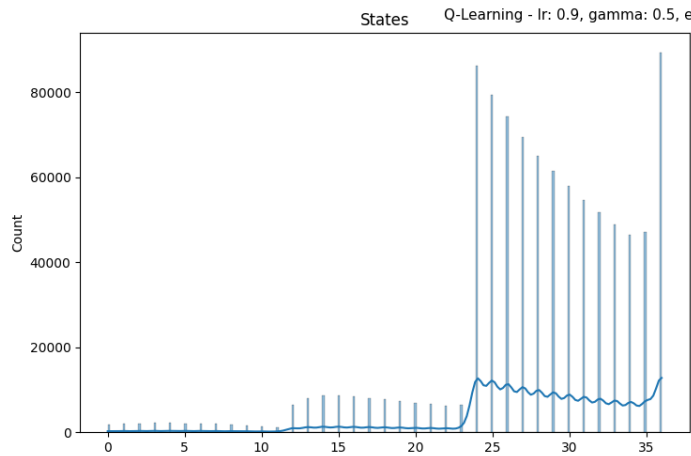
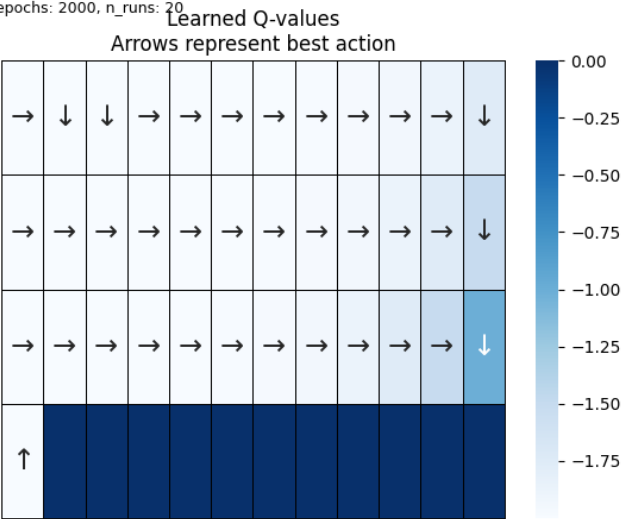
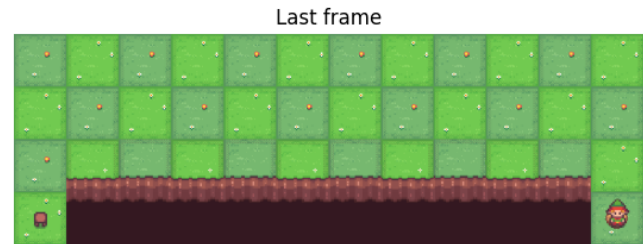
Cumulative Rewards

[illegible]

	Episodes	Rewards	Steps	cum_rewards	map_size
0	0	-1498.0	310.0	-1498.0	4x12
1	1	-893.0	299.0	-2391.0	4x12
2	2	-52.0	52.0	-2443.0	4x12
3	3	-48.0	48.0	-2491.0	4x12
4	4	-45.0	45.0	-2536.0	4x12
...
39995	1995	-13.0	13.0	-102033.0	4x12
39996	1996	-13.0	13.0	-102046.0	4x12
39997	1997	-15.0	15.0	-102061.0	4x12
39998	1998	-13.0	13.0	-102074.0	4x12
39999	1999	-15.0	15.0	-102089.0	4x12

Eksperyment 3:

Q-Learning - lr: 0.9, gamma: 0.5, epsilon: 0.2, epochs: 2000, n_runs: 20



Qtable

[-1.99992895	-1.99987792	-1.99987793	-1.99991949]
[-1.99985888	-1.99975586	-1.99975586	-1.99990275]
[-1.99972761	-1.99951172	-1.99951172	-1.99984119]
[-1.99945005	-1.99902344	-1.99902344	-1.99963212]
[-1.99896332	-1.99804687	-1.99804688	-1.99932922]
[-1.99798912	-1.99609375	-1.99609375	-1.99890351]
[-1.99571417	-1.9921875	-1.9921875	-1.99799439]
[-1.99145251	-1.984375	-1.984375	-1.99560304]
[-1.98279157	-1.96875	-1.96875	-1.99200894]
[-1.9666849	-1.9375	-1.9375	-1.98067256]
[-1.93282019	-1.875	-1.875	-1.96640919]
[-1.86616634	-1.86101123	-1.75	-1.92456297]
[-1.99993879	-1.99975586	-1.99975586	-1.99987793]
[-1.99987793	-1.99951172	-1.99951172	-1.99987793]
[-1.99975586	-1.99902344	-1.99902344	-1.99975586]
[-1.99951172	-1.99804688	-1.99804688	-1.99951172]
[-1.99902344	-1.99609375	-1.99609375	-1.99902344]
[-1.99804687	-1.9921875	-1.9921875	-1.99804687]
[-1.99609375	-1.984375	-1.984375	-1.99609375]
[-1.9921875	-1.96875	-1.96875	-1.9921875]
[-1.984375	-1.9375	-1.9375	-1.984375]
[-1.96875	-1.875	-1.875	-1.96875]
[-1.9375	-1.75	-1.75	-1.93749996]
[-1.87499896	-1.75	-1.5	-1.875]
[-1.99987793	-1.99951172	-1.99987793	-1.99975586]
[-1.99975586	-1.99902344	-100.99987793	-1.99975586]
[-1.99951172	-1.99804688	-100.99987793	-1.99951172]
[-1.99902344	-1.99609375	-100.99987793	-1.99902344]
[-1.99804688	-1.9921875	-100.99987793	-1.99804688]
[-1.99609375	-1.984375	-100.99987793	-1.99609375]
[-1.9921875	-1.96875	-100.99987793	-1.9921875]
[-1.984375	-1.9375	-100.99987793	-1.984375]
[-1.96875	-1.875	-100.99987793	-1.96875]
[-1.9375	-1.75	-100.99987793	-1.9375]
[-1.875	-1.5	-100.99987793	-1.875]
[-1.75	-1.5	-1.	-1.75]
[-1.99975586	-100.99987793	-1.99987793	-1.99987793]
[0.	0.	0.	0.]
[0.	0.	0.	0.]
[0.	0.	0.	0.]
[0.	0.	0.	0.]
[0.	0.	0.	0.]
[0.	0.	0.	0.]
[0.	0.	0.	0.]
[0.	0.	0.	0.]
[0.	0.	0.	0.]
[0.	0.	0.	0.]
[0.	0.	0.	0.]
[0.	0.	0.	0.]

Cumulative Rewards

	Episodes	Rewards	Steps	cum_rewards	map_size
0	0	-1807.0	520.0	-1807.0	4x12
1	1	-27.0	27.0	-1834.0	4x12
2	2	-215.0	116.0	-2049.0	4x12
3	3	-179.0	80.0	-2228.0	4x12
4	4	-48.0	48.0	-2276.0	4x12
...
39995	1995	-130.0	31.0	-221996.0	4x12
39996	1996	-229.0	31.0	-222225.0	4x12
39997	1997	-13.0	13.0	-222238.0	4x12
39998	1998	-17.0	17.0	-222255.0	4x12
39999	1999	-15.0	15.0	-222270.0	4x12

Obserwacje i wnioski

1. Stan 35 – w którym odnotować można najwięcej wykonanych kroków/akcji podczas treningu to stan początkowy, w którym gracz/uczeń pojawia się zaczynając epizod, natomiast stan 24, który ma równie wysoką liczbę odwiedzin to jedyny stan do którego można przejść ze stanu początkowego (nie wpadając z klifu), są to dwa kroki które uczeń musi zrobić by kontynuować naukę. Najczęściej odwiedzane stany są w przedziale 24-36, czyli stany bliskie klifu – reprezentujące najlepszą ścieżkę. Pokazuje to że uczeń w trakcie treningu uczył się najlepszej ścieżki i eksploatował jej możliwości z niewielką eksploracją pozostałych miejsc.
2. Również najczęściej podejmowane przez ucznia akcje to pójście w prawo, co jest zgodne z wizualizacją na heatmap (średnio w większości stanów z oceny środowiska wychodzi, że największą wygraną (tu: jako najmniejszą stratę) zapewni pójście w lewo)
3. Współczynnik gamma (w kodzie: `discount_factor`) wpływa na to jak istotny w ocenie wybranej akcji/przejścia ze stanu $s_0 \rightarrow s_1$ jest nagroda jaką można dostać robiąc krok $s_1 \rightarrow s_2$. W skrócie powinna być $\gamma \neq 1$ – jeśli gamma byłaby równa 1 to algorytm nie brałby w ogóle pod uwagę przy wyborze przejścia przyszłych możliwych nagród, z kolei nawet przy dosyć niskim współczynniku gamma (jak $\gamma=0.5$) a odpowiednio dobranych pozostałych hiperparametrach algorytm w niewielkim stopniu modyfikuje q-table (Eksperyment 3) a jest w stanie nauczyć się najlepszej strategii.
4. Algorytm Q-learning jest dobrym wyborem w przypadku rozwiązywania niezbyt skomplikowanych problemów w prostych środowiskach – np. w takim jak Cliff Walking, dodatkowo w małych środowiskach nauka jest bardzo szybka. Uczeń ‘zapamiętuje’ najlepszą ścieżkę (jako wartości qtable) i podczas każdego epizodu/epoki modyfikuje ją na podstawie nowych ‘obserwacji’ w ten sposób wyznaczając optymalne rozwiązanie. Natomiast algorytm nie dałby dobrych rezultatów w bardziej złożonych problemach – gdzie np. środowisko modyfikuje się w każdej epoce.