

WSI Uczenie ze wzmocnieniem sprawozdanie

Małgorzata Grzanka

26.05.2024

Spis treści

1	Wstęp	1
2	Zasada działania Q-learningu	2
2.1	Ogólny koncept	2
2.2	Uczenie się	2
2.3	Wybór akcji	2
3	Hiperparametry	3
4	Wpływ liczby epizodów	3
4.1	Wyniki	3
4.2	Wnioski	4
5	Wpływ wartości dyskontu	5
5.1	Wyniki	5
5.2	Wnioski	5
6	Wpływ wartości learning rate	6
6.1	Wyniki	6
6.2	Wnioski	7
7	Wpływ wartości T	8
7.1	Wyniki	8
7.2	Wnioski	8

1 Wstęp

Celem ćwiczenia było zaimplementowanie algorytmu Q-learning. Następnie, algorytm należało przetestować, rozwiązując za jego pomocą problem Taxi, dostępny w pakiecie gymnasium ([https:// gymnasium.farama.org/environments/toy_text/taxi/](https://gymnasium.farama.org/environments/toy_text/taxi/)).

2 Zasada działania Q-learningu

2.1 Ogólny koncept

Q-learning polega na uczeniu agenta najoptymalniejszych zachowań w jakimś środowisku.

Najważniejszym elementem w procesie uczenia jest Q-function. Q-function w jakimś stanie s_t po akcji a_t to suma nagród, które gracz osiągnie, robiąc akcje a_t i dalej wykonując ruchy zgodnie z polityką π . Zatem agent powinien w większości przypadków wybierać akcje, dla której, przy danym stanie środowiska, Q-function zwróci największą wartość.

$$Q^\pi(s_t, a_t) = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \gamma^3 r_{t+3} + \dots$$

W Q-learningu, jako że wartości Q-function są na początku inicjowane losowymi wartościami, możemy obliczać wartości Q-function za pomocą jej wartości dla kolejnego stanu (wybieramy największą wartość Q-function dla kolejnego ruchu):

$$Q^\pi(s_t, a_t) = r_t + \gamma \cdot \max_{a_{t+1}} (Q^\pi(s_{t+1}, a_i))$$

Dla mniej skomplikowanego środowiska, jakim jest Taxi problem, wszystkie wartości, które może przyjąć Q-function, można zapisać w postaci tabeli o wymiarach $|S| \times |A|$, gdzie $|S|$ to liczba wszystkich możliwych stanów w grze, a $|A|$ to liczba możliwych akcji, które można w tej grze podjąć.

2.2 Uczenie się

Proces uczenia polega na wyznaczeniu przez algorytm odpowiednich wartości Q-function poprzez wielokrotne ruchy agenta w środowisku i otrzymywanie od tego środowiska sprzężenia zwrotnego w postaci nagrody. Algorytm dąży do tego, aby różnica pomiędzy wartością Q-function otrzymanej z tabeli była najbliższa możliwa wartości wyliczonej za pomocą wzoru na Q-function.

Stosuje się gradient descend do minimalizacji tego błędu:

$$Q^\pi(s_t, a_t) = Q^\pi(s_t, a_t) + \alpha \cdot (r_t + \gamma \cdot \max_{a_{t+1}} (Q^\pi(s_{t+1}, a_i)) - Q^\pi(s_t, a_t))$$

2.3 Wybór akcji

Wybór jakiejś akcji przez agenta odbywa się z prawdopodobieństwami wyliczonymi na podstawie strategii boltzmannowskiej:

$$P(a|s) = \frac{e^{\frac{Q(s,a)}{T}}}{\sum_{a'} e^{\frac{Q(s,a')}{T}}}$$

Parametr T (temperatura) kontroluje balans między eksploracją a eksploatacją w strategii Boltzmannowskiej. Wysoka temperatura sprzyja eksploracji (losowemu wybieraniu akcji), podczas gdy niska temperatura sprzyja eksploatacji (wybieraniu najlepszej znanej akcji).

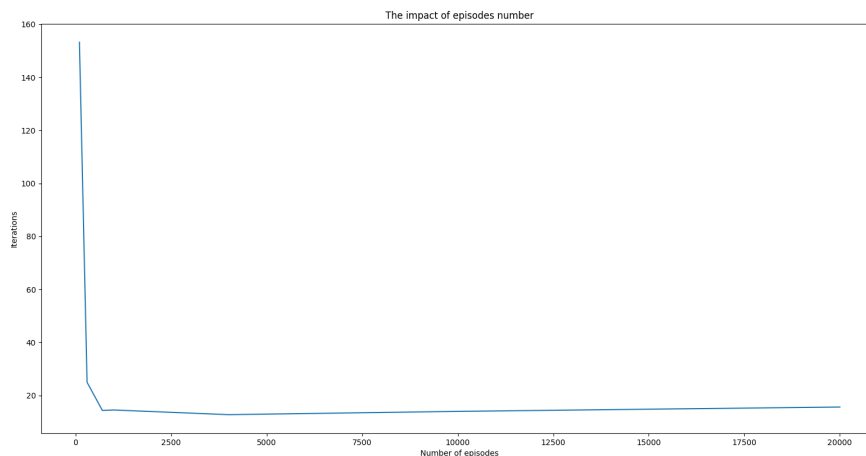
3 Hiperparametry

Mój model algorytmu Q-learningu przyjmuje następujące hiperparametry:

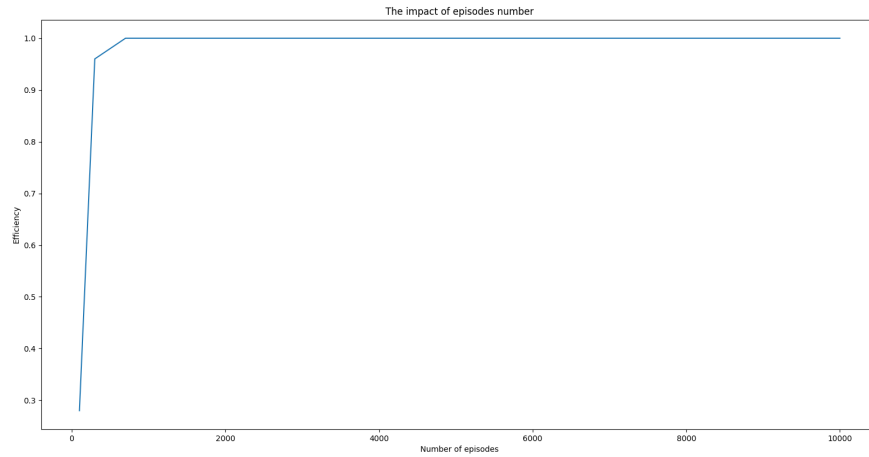
- dyskonto - jak bardzo przyszłe nagrody są brane pod uwagę przy podejmowaniu bieżących decyzji.
- T - temperatura, kontroluje balans między eksploracją a eksploatacją w strategii Boltzmannowskiej.
- epizody - ile symulacji treningowych algorytm powinien rozegrać, aby się nauczyć.
- learning rate - współczynnik dla minimalizacji błędu Q -function (dla algorytmu gradient descent).

4 Wpływ liczby epizodów

4.1 Wyniki



Rysunek 1: Wpływ liczby epizodów na liczbę iteracji podczas symulacji



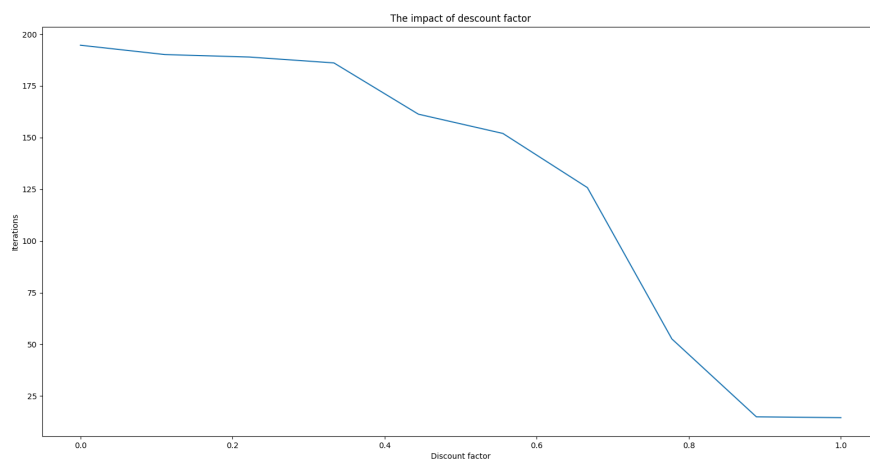
Rysunek 2: Wpływ liczby epizodów na skuteczność agenta podczas symulacji

4.2 Wnioski

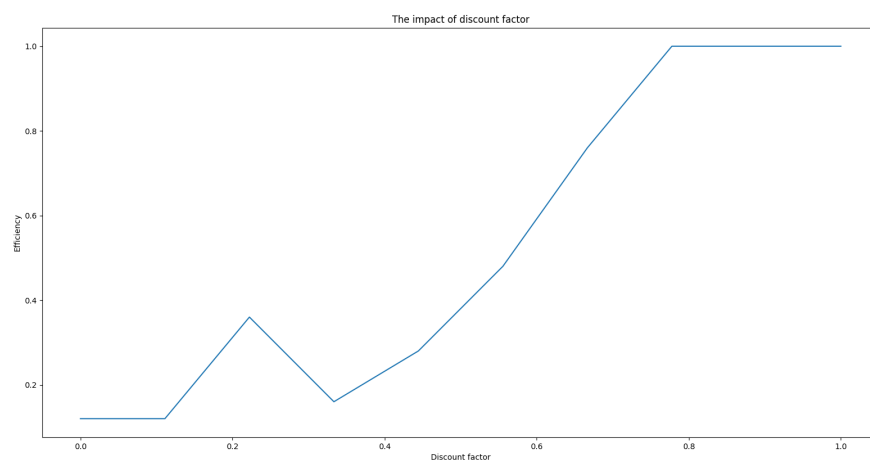
Wraz z wzrostem epizodów, wartości przyjmowane przez Q-function są bardziej dopasowane dla danego środowiska (więcej prób podczas treningu agenta). Dlatego, dla małej ilości epizodów, wytrenowanemu agentowi ukończenie symulacji zajmuje bardzo dużo czasu (ciagle wykonuje niepoprawne ruchy), a jego skuteczność jest bardzo słaba. Liczba iteracji potrzebna do zakończenia symulacji przez wytrenowanego agenta spada gwałtownie już przy około 400 próbach, ale swoje absolutne minimum (około 12 iteracji) osiąga dla liczby epizodów równej około 4000. Jeśli chodzi o skuteczność, jest ona równa 100% dla liczby epizodów większej niż 1000.

5 Wpływ wartości dyskontu

5.1 Wyniki



Rysunek 3: Wpływ wartości dyskontu na liczbę iteracji podczas symulacji



Rysunek 4: Wpływ wartości dyskontu na skuteczność agenta podczas symulacji

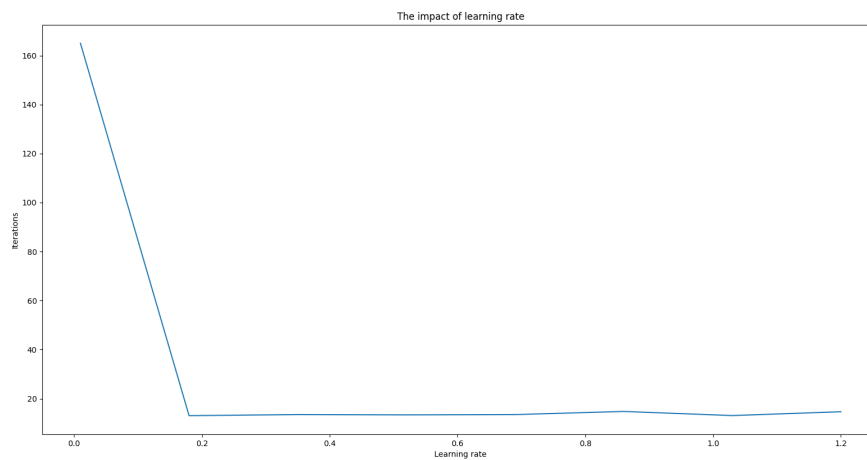
5.2 Wnioski

Współczynnik dyskontowy określa, jak bardzo przyszłe nagrody są uwzględniane w porównaniu do bieżących nagród i przyjmuje wartości od 0 do 1.

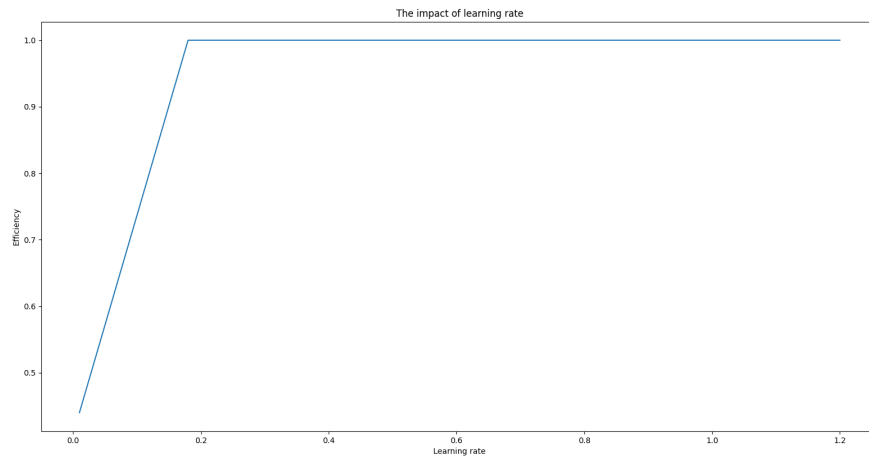
Dla rozważanego problemu, lepsza skuteczność i szybsze dotarcia do końca symulacji osiągane są dla wyższych wartości dyskontu. Małe wartości dyskontu sprzyjają dużym, natychmiastowym nagrodom, a większe prowadzi do bardziej optymalnych strategii w dłuższym horyzoncie czasowym. Jednak może to również sprawić, że agent będzie bardziej wrażliwy na zmienność i stochastyczność środowiska, co może utrudniać szybkie zdobycie nagród. Dla tej symulacji, lepsze okazały się te większe wartości współczynnika dyskontu.

6 Wpływ wartości learning rate

6.1 Wyniki



Rysunek 5: Wpływ wartości learning rate na liczbę iteracji podczas symulacji



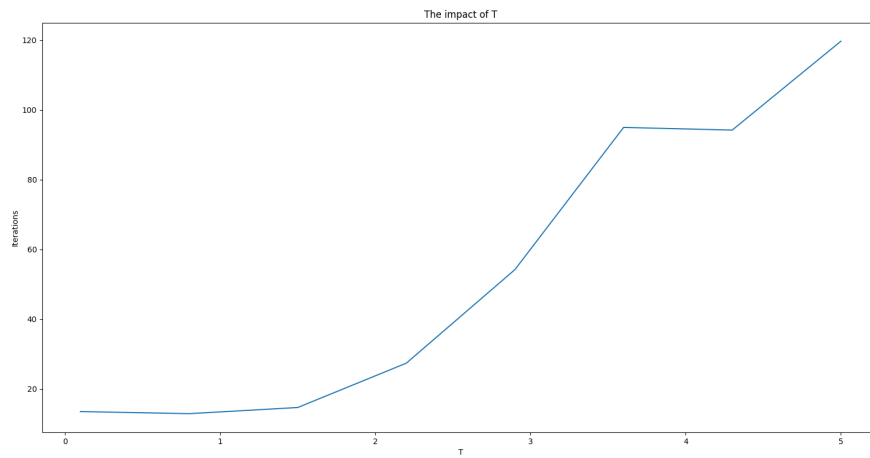
Rysunek 6: Wpływ wartości learning rate na skuteczność agenta podczas symulacji

6.2 Wnioski

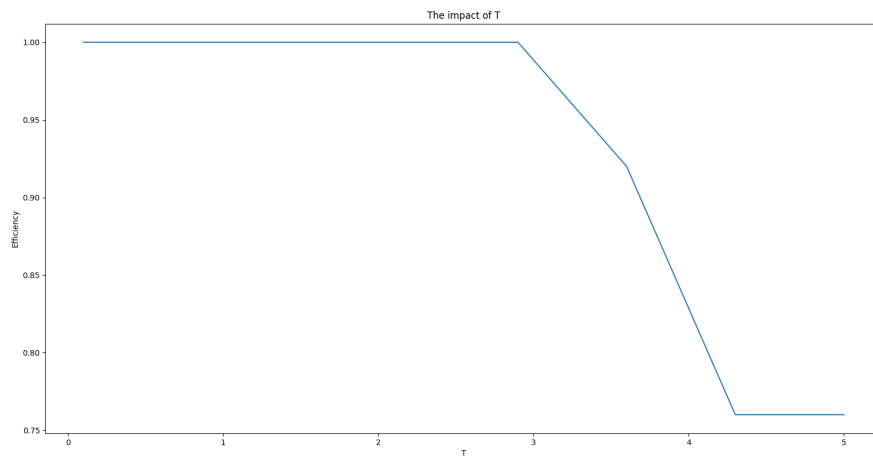
Im większe learning rate, tym szybciej agent uczy się właściwych ruchów na symulacji. Dla bardzo małych learning rate, liczba epizodów równa 3000 nie była wystarczająca dla właściwego nauczenia agenta - liczba iteracji była bardzo duża, a skuteczność bardzo mała. Jednakże, dla bardzo dużych learning rate, wartości Q-function stają się za duże dla poprawnego policzenia prawdopodobieństwa strategii Boltzman (overflow w podnoszeniu e to potęgi). Idealne wartości dla tego problemu są wartości z przedziału 0.4 - 1.

7 Wpływ wartości T

7.1 Wyniki



Rysunek 7: Wpływ wartości T na liczbę iteracji podczas symulacji



Rysunek 8: Wpływ wartości T na skuteczność agenta podczas symulacji

7.2 Wnioski

Im większe T, tym wyliczane prawdopodobieństwa mniej zależą od wartości Q-function dla możliwych akcji. Oznacza to, że prawdopodobieństwa wykonania

każdej akcji w danym stanie są do siebie zbliżone. Oznacza to, że agent jest nastawiony bardziej na eksplorację nowych rozwiązań. Z kolei dla małych wartości T , agent skupia się bardziej na eksploatacji znalezionych wcześniej wartości Q -function. W przypadku problemu Taxi, bardziej efektywne są mniejsze wartości T (z przedziału 0 - 1).