

Sprawozdanie WSI - ćwiczenie 6 Q Learning

Zaimplementuj algorytm Q-learning. Następnie, wykorzystując proste środowisko (np. Taxi-v3), zbadaj wpływ hiperparametrów na działanie algorytmu (np. wpływ strategii eksploracji, współczynnik uczenia).

Pamiętaj, że implementacja musi być wykonana samodzielnie. Brak zrozumienia dostarczonego kodu rozwiązania równoważny jest plagiatowi!

Rozwiązanie:

Zadanie zostało rozwiązane przy pomocy algorytmu q learning z zastosowaniem dwóch strategii eksploracji. Strategia ϵ -zachłanna i strategia boltzmannna.

Strategia z użyciem epsilon , polega na wybraniu z pewnym prawdopodobieństwem $\epsilon > 0$ dowolną akcję losowo według rozkładu równomiernego, a z prawdopodobieństwem $1-\epsilon$ wybiera się akcję zachłanną (jeśli jest ich wiele, to także losowo). Formalnie można to zapisać następująco:

$$\pi(x, a^*) = \begin{cases} \frac{1-\epsilon}{|\text{Arg max}_a Q(x, a)|} + \frac{\epsilon}{|A|} & \text{jeśli } a^* \in \text{Arg max}_a Q(x, a), \\ \frac{\epsilon}{|A|} & \text{w przeciwnym przypadku.} \end{cases}$$

Strategia z zachłannym epsilonm posiada wadę, że prawdopodobieństwo losowego zachowania się ucznia nie zależy od tego czego zdołał się nauczyć. Jednym ze sposobów pozbycia się tego problemu jest zastosowanie strategii boltzmannna. Która jest opisana wzorem:

$$\pi(x, a^*) = \frac{\exp(Q(x, a^*)/T)}{\sum_a \exp(Q(x, a)/T)},$$

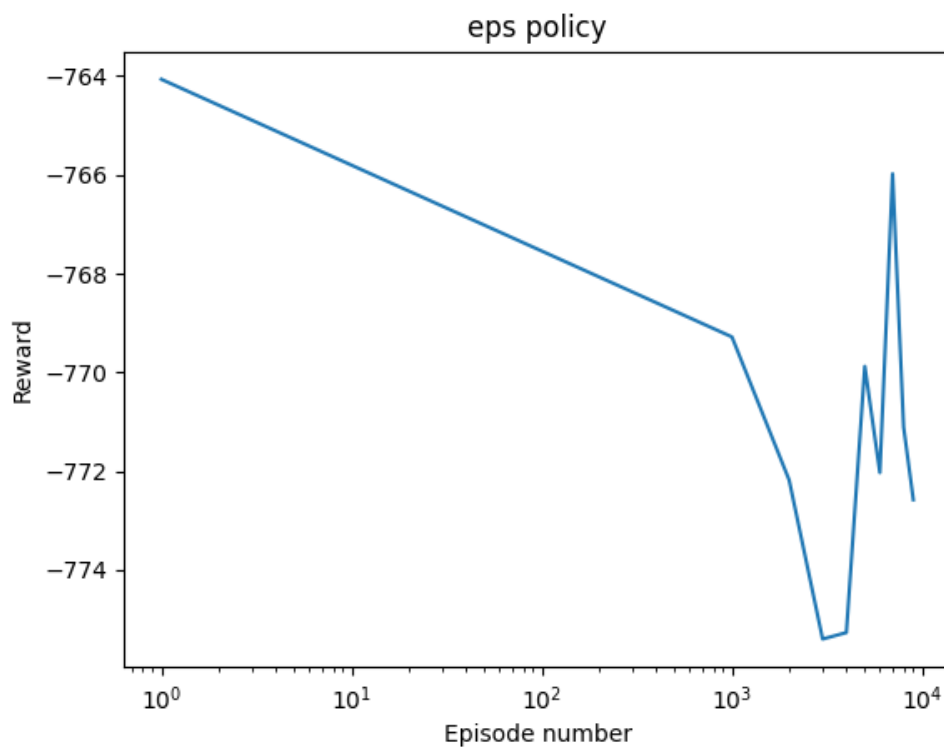
Współczynnik T (temperatura) > 0 reguluje stopień losowości. Przy tej strategii prawdopodobieństwo wyboru akcji niezachłannej jest tym mniejsze, im bardziej akcja zachłanna ma większą Q -wartość od pozostałych.

Testy:

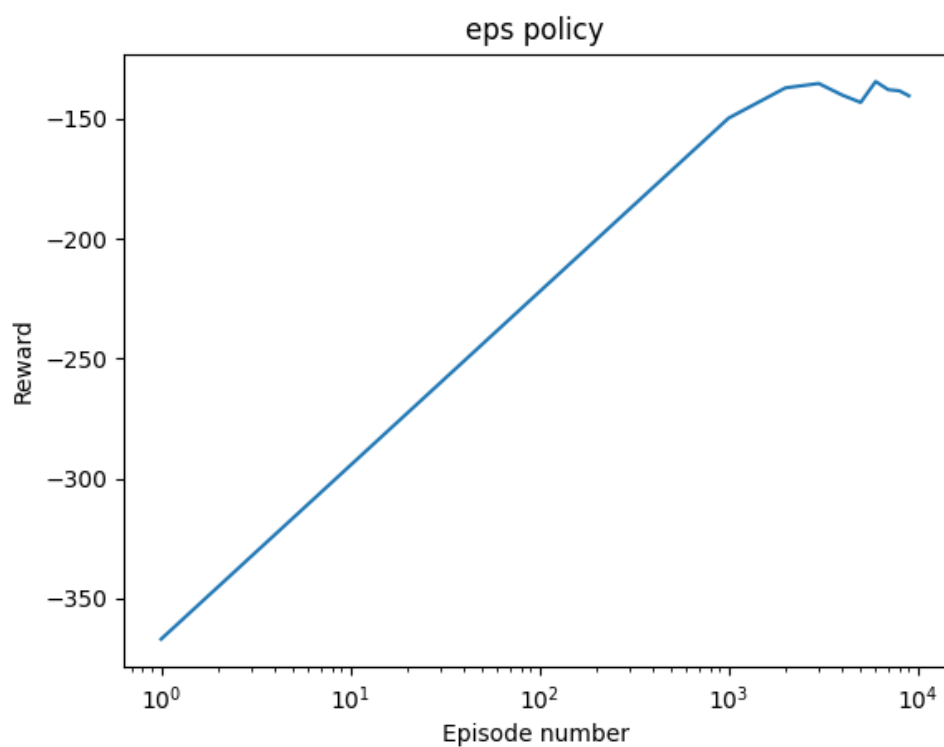
Strategia zachłannego epsilon:

Badanie zachowania algorytmu przy zmianie eps. Pozostałe parametry (discount factor = 1, learning rate = 0.1, episodes = 10000, steps = 1000)

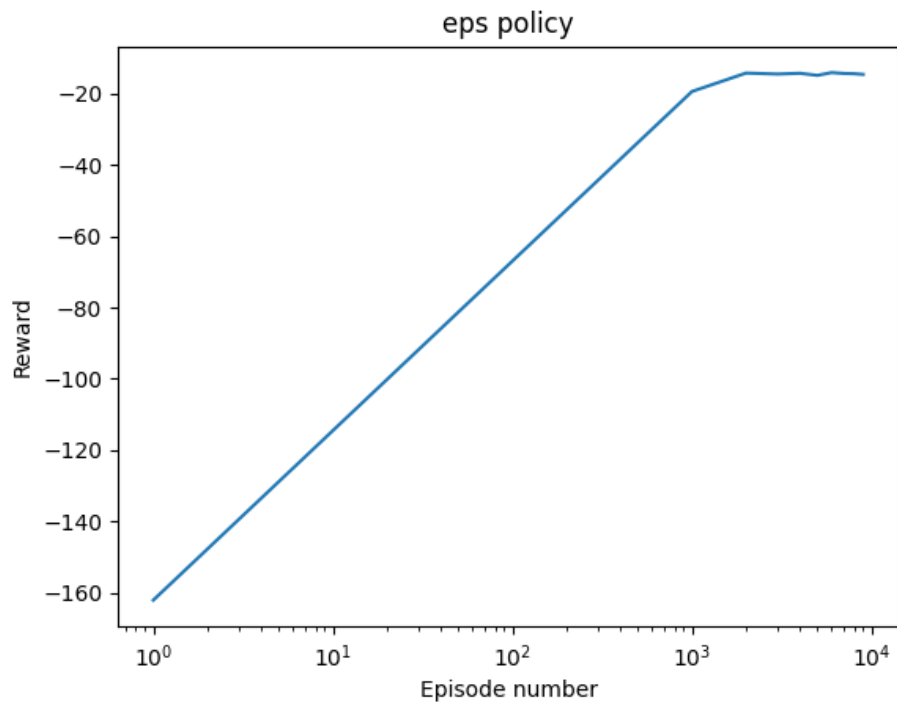
Przy epsilon = 1 zgodnie z definicją strategii zauważymy 100% losowość:



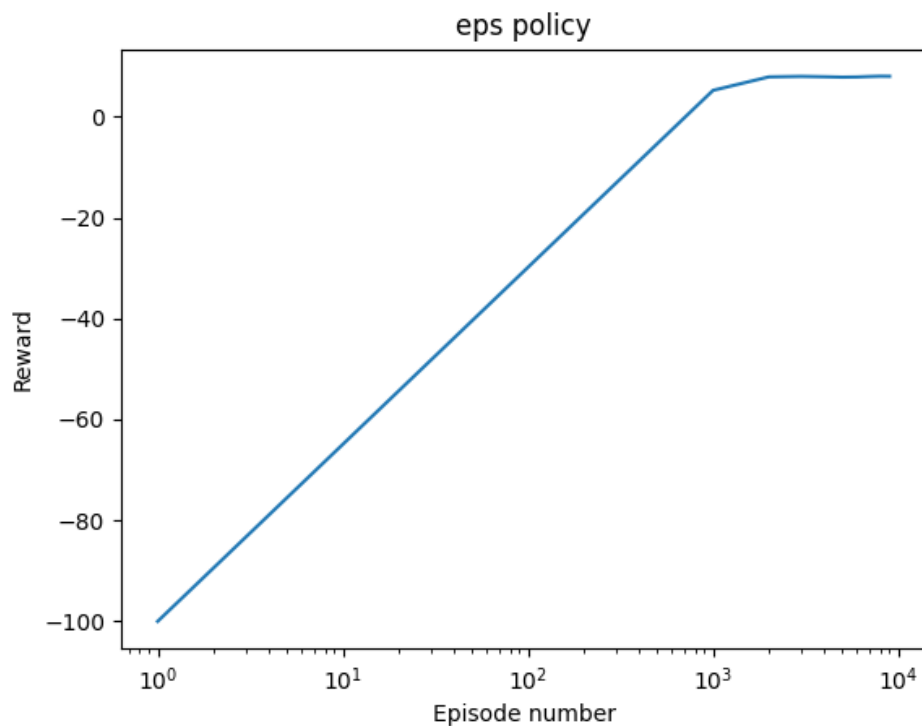
Przy stopniowym zmniejszaniu epsilon dochodzimy do coraz lepszych wyników i jednocześnie wyższych nagród. Poniżej wykres dla $\epsilon=0.7$



Dla $\epsilon=0.3$

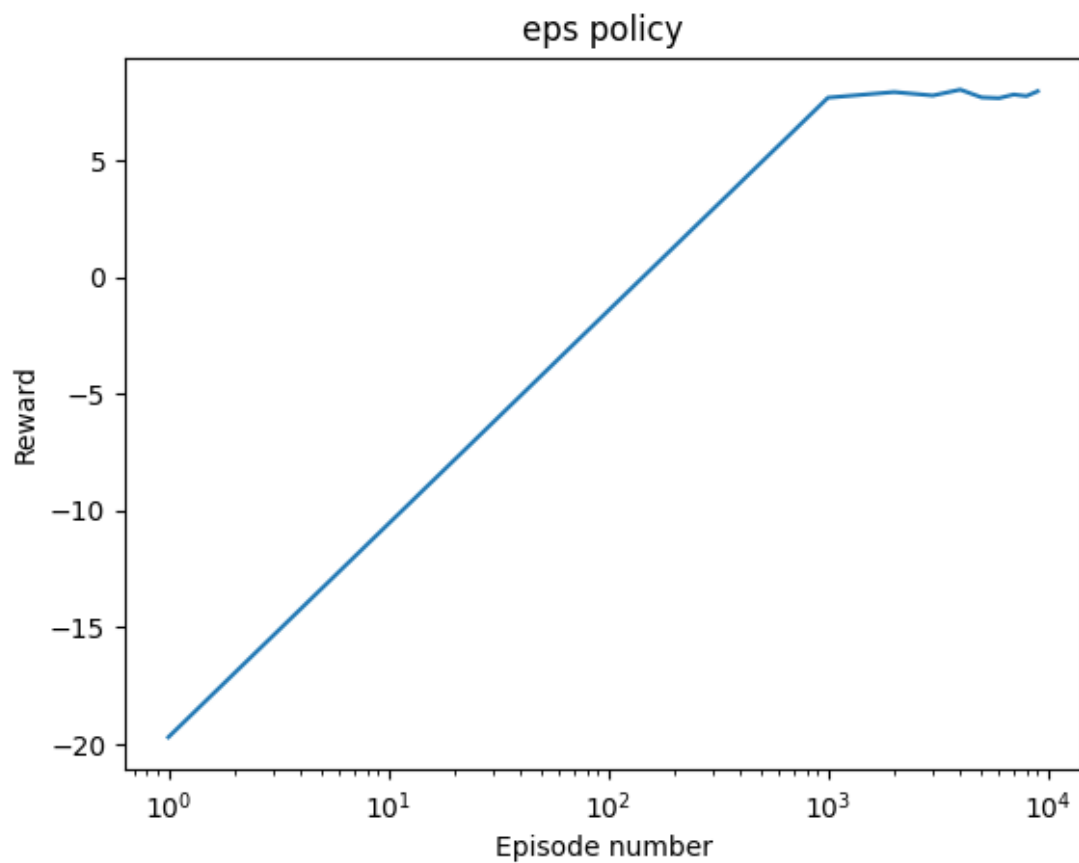


Dla epsilon równego 0 otrzymujemy najlepsze wyniki



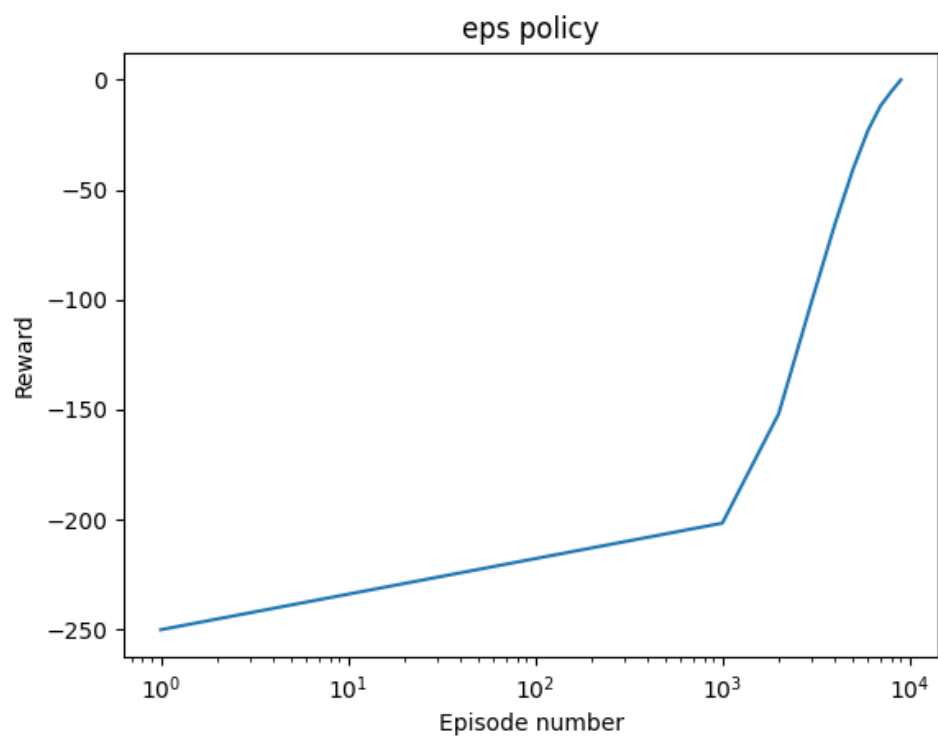
Badanie wpływu zmiany parametru learning rate:

Podczas gwałtownej zmiany learning rate do wartości 0.9 możemy zaobserwować wzrost maksymalnej wartości nagrody. ($\epsilon = 0$)

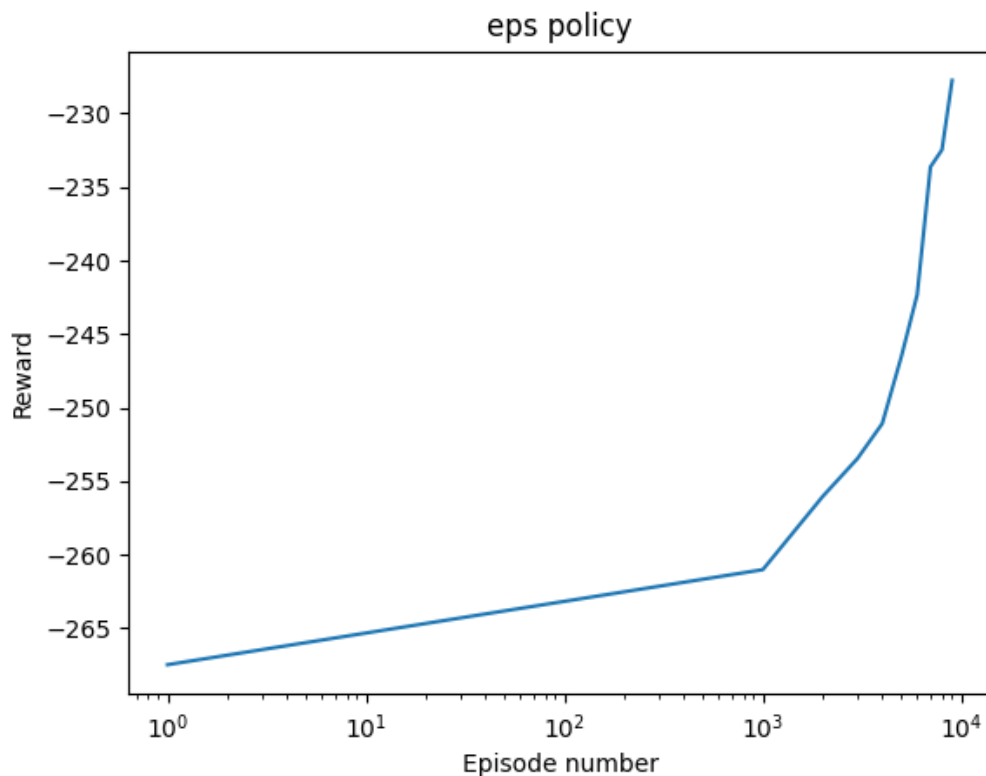


Natomiast dla dużo mniejszych wartości learning rate, algorytm potrzebowałby więcej iteracji aby otrzymał akceptowalne wyniki.

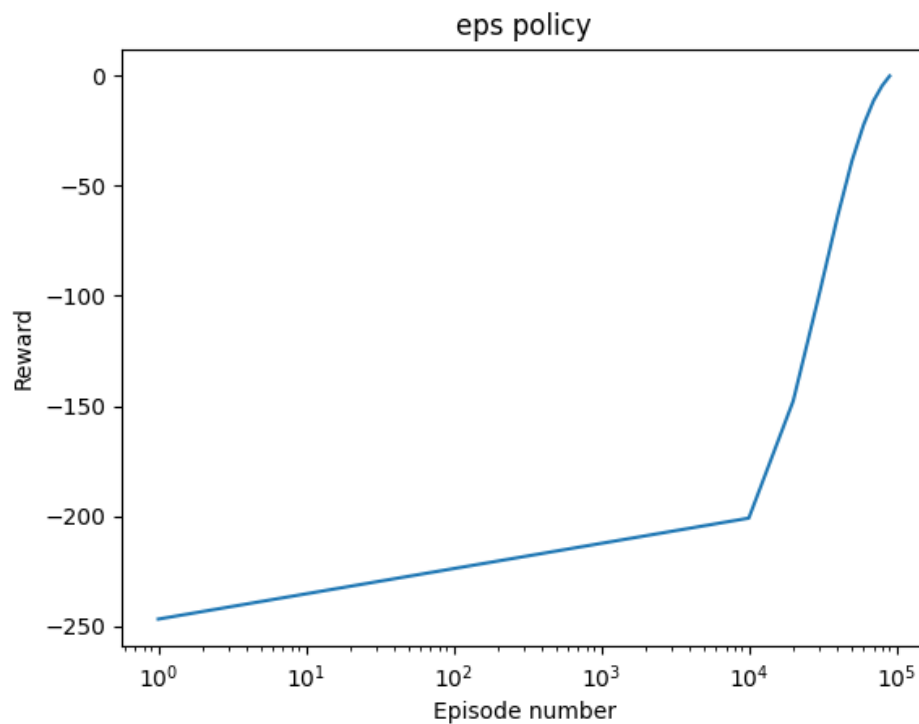
Dla learning rate = 0.01



W przypadku zmniejszenia wartości do 0.001



można zauważyć że w tym przypadku algorytm miał przyjętą za małą ilość epizodów jak na tak mały learning rate. W momencie zwiększenia ilości epizodów algorytm mimo bardzo małej wartości learning rate dochodzi do średniej wartości nagród w okolicy 0.



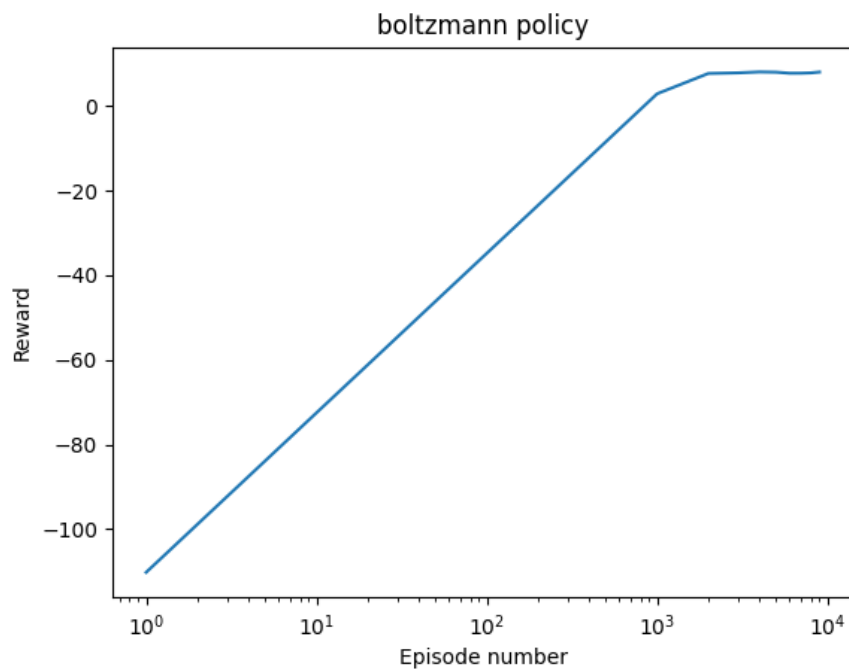
Wnioski: metoda zachłannego epsilon przy dobraniu odpowiednich parametrów spełnia swoje zadanie. Zmiana epsilon działa zgodnie z teorią. Manipulacja parametrami dot. learning rate i ilości epizodów ma znaczący wpływ na wynik końcowy algorytmu. W

przypadku dobrania mniejszego learning rate'u należy zwiększać liczbę epizodów jeżeli chcemy otrzymać podobne wyniki.

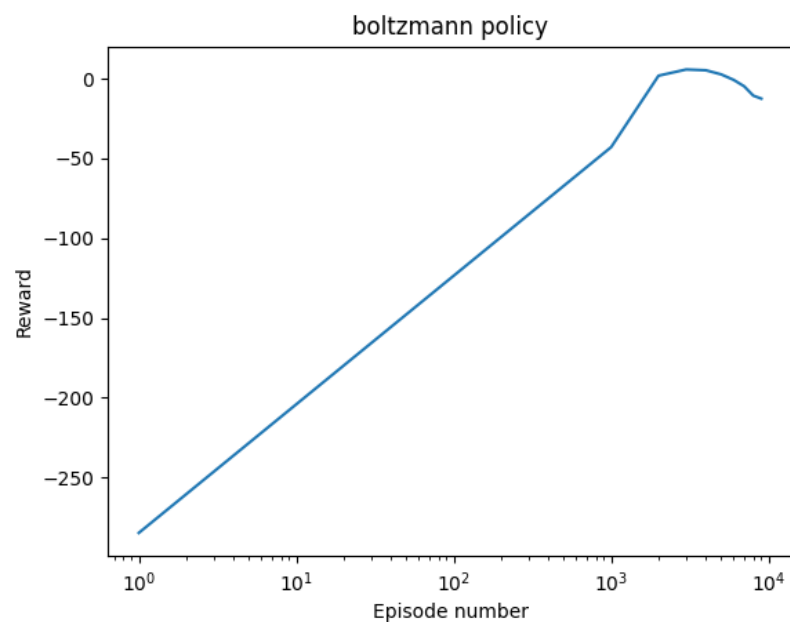
Strategia Boltzmannna

Wpływ parametru temperatury na algorytm:

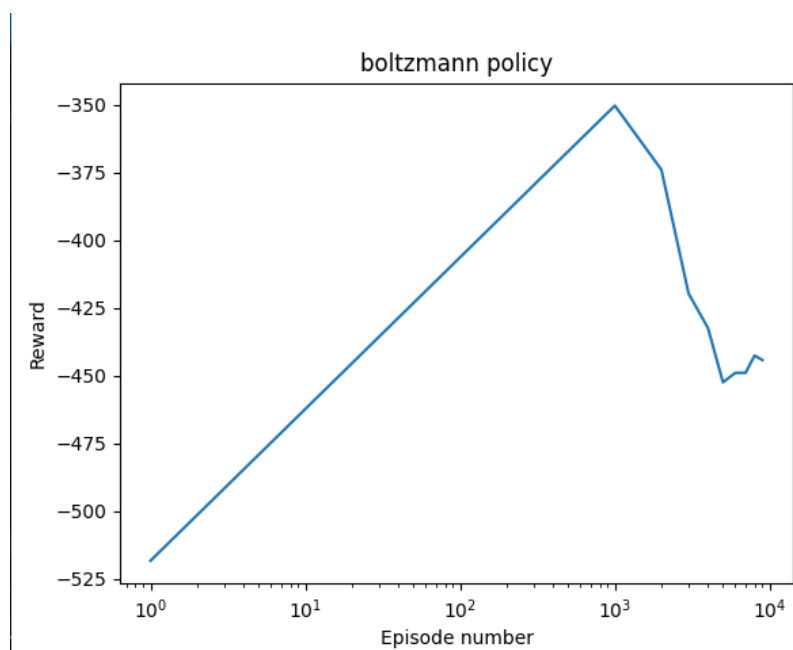
Dla $T = 0.1$



Dla $T = 3$

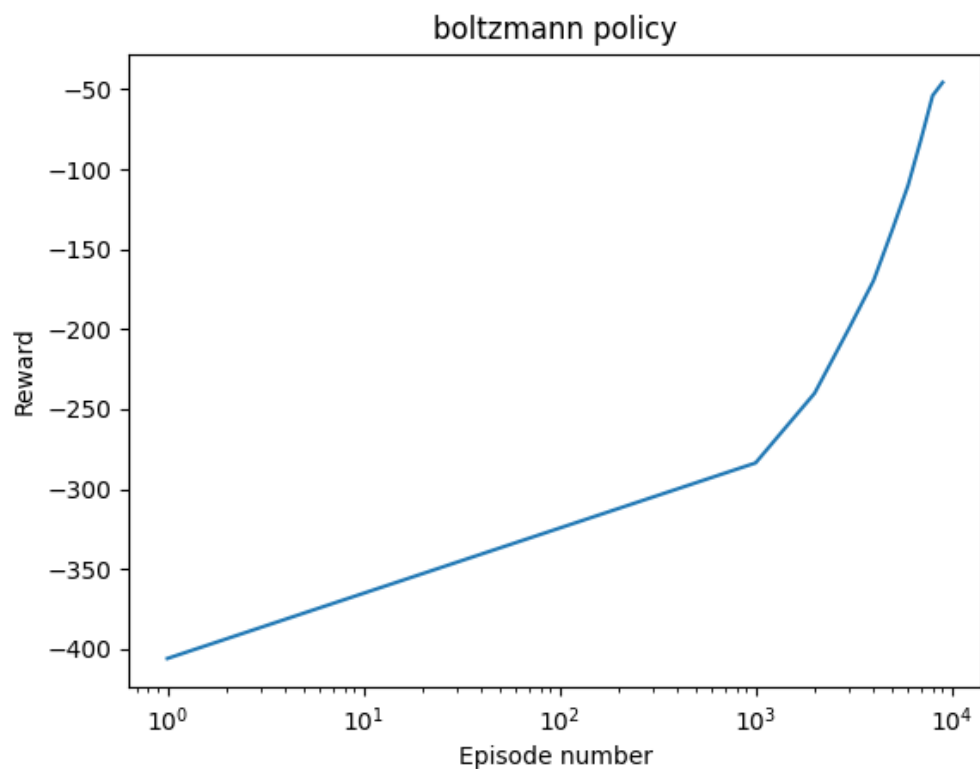


Dla $T = 10$



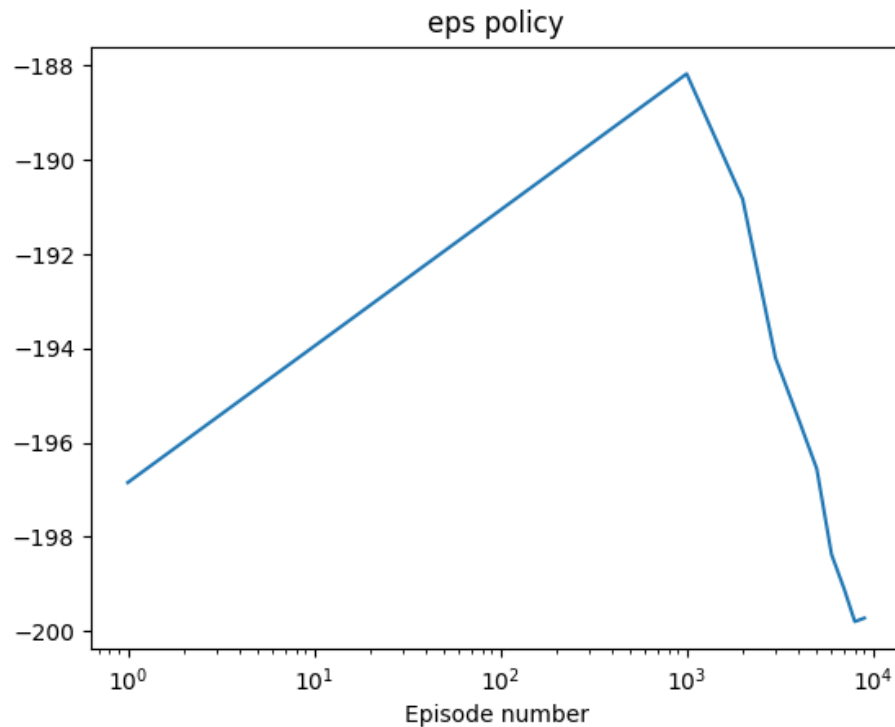
Zwiększając parametr temperatury algorytm zaczyna działać z większą losowością co jest zgodne z oczekiwaniami.

Strategia boltzmana zachowuje się podobno do strategii zachłannego epsilon jeżeli chodzi o wpływ learning rate na wykres nagrody od ilości epizodów. Dla małego learning rate = 0.01

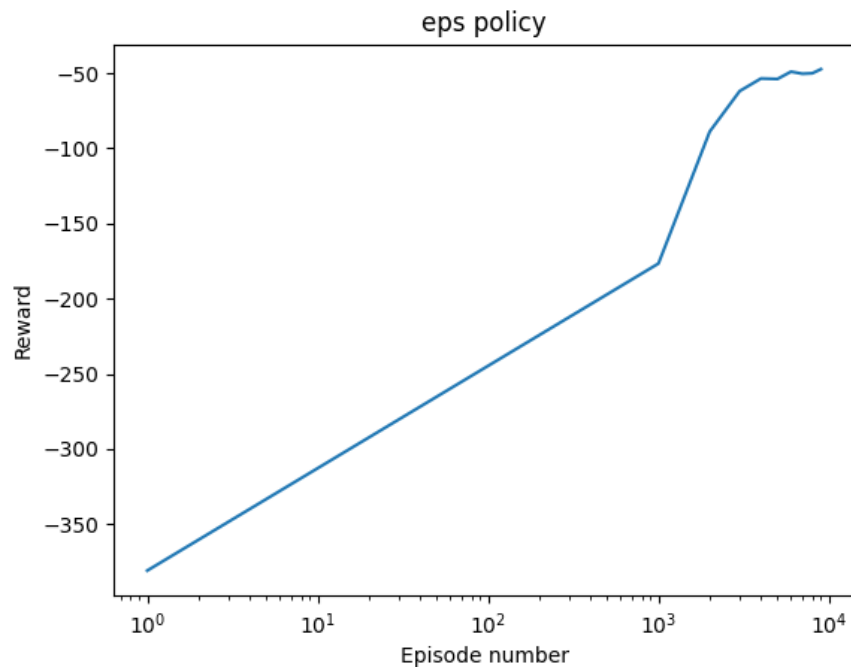


Wpływ wartości discount factor:

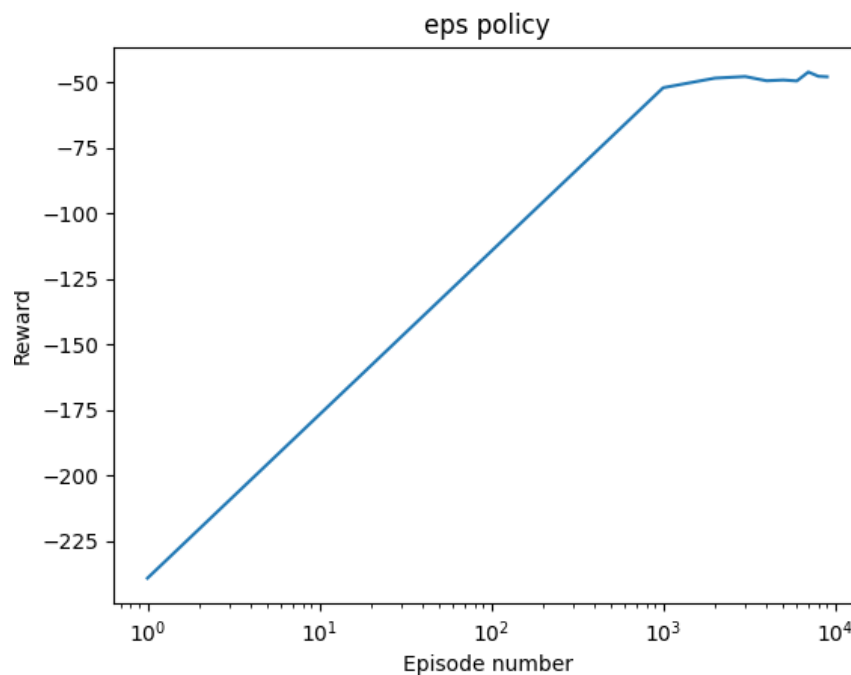
W przypadku ustawienia wartości discount factor na 0, algorytm nie będzie patrzył 'w przód' i będzie zwracał uwagę jedynie na natychmiastowe nagrody co spowoduje brak otrzymania prawidłowego rozwiązania zadania.



Natomiast przy tym parametrze ustawionym na wartość 0.3 przy epsilon = 0.5



discount factor = 1 przy eps = 0.5



Wnioski:

Algorytm działa zgodnie z założeniami. Przy doborze odpowiednich parametrów zarówno strategia boltzmannna jak i zachłannego epsilon radzi sobie z uczeniem i zrealizowaniem zadania taxi. Najbardziej znaczącymi parametrami jest epsilon i temperatura, od których zależy skuteczność/prawdopodobieństwo wybrania losowego. Discount factor również nie może być zbyt mały ponieważ algorytm będzie brał pod uwagę jedynie natychmiastowe nagrody