

Sprawozdanie Laboratorium 6 Uczenie ze wzmocnieniem

Piotr Gierżatowicz-Sierpień 331376

Cel zadania

Należało zaimplementować algorytm qlearning, który miał wskazać optymalną politykę dla agenta zarządzającego stanem magazynowym sklepu. Agent otrzymuje nagrodę = 1 przy akcji sprzedaży gdy stan > 0, oraz nagrodę = 100 przy akcji kupienia, gdy stan = 9. Należało również przeprowadzić wpływ parametrów algorytmu na znajdowanie optymalnej polityki.

Implementacja

Parametry wykorzystane w badaniu zachowania algorytmu:

Alpha = 0.0005, 0.001, 0.005, 0.01, 0.02, 0.05, 0.1, 0.15, 0.2, 0.3 <- learning_rate

Gamma = 1.0

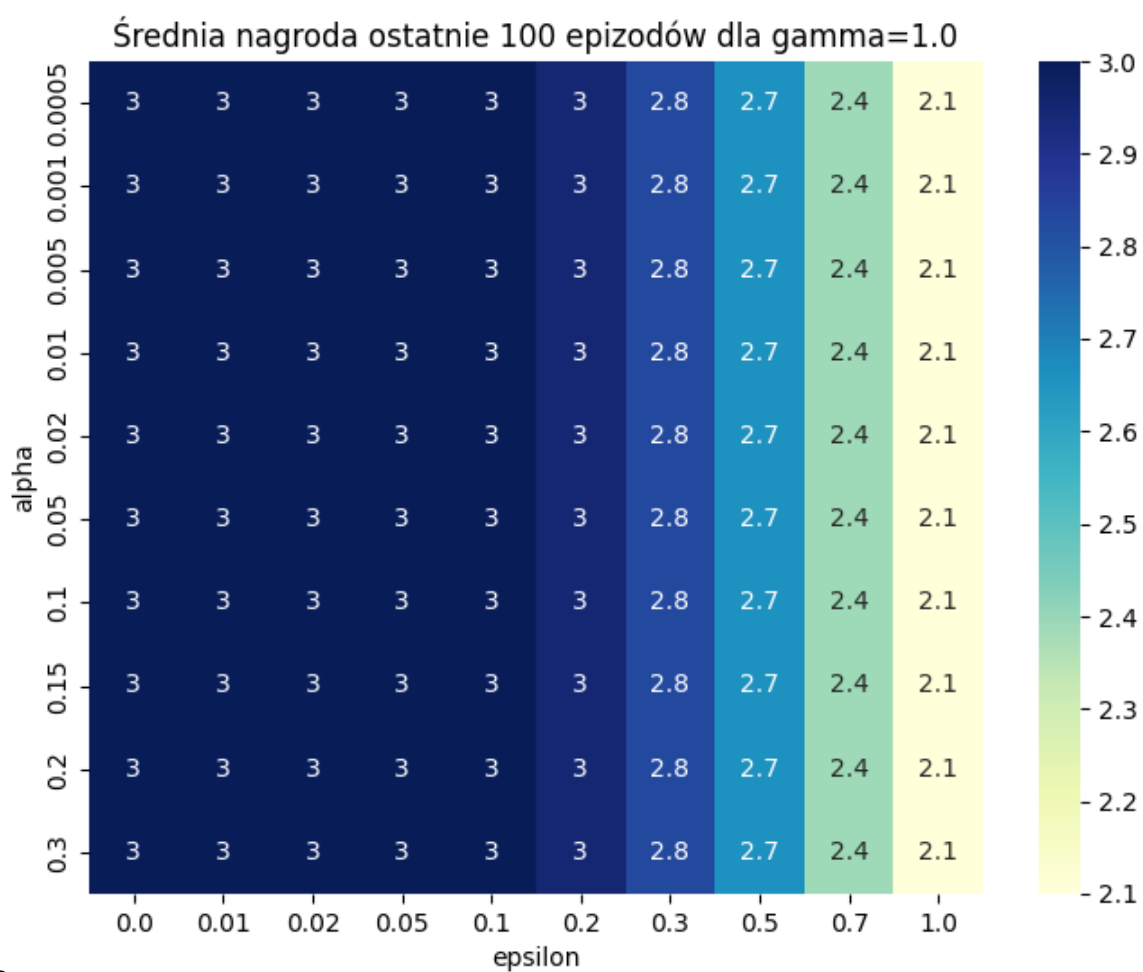
Epsilon = 0.0, 0.01, 0.02, 0.05, 0.1, 0.2, 0.3, 0.5, 0.7, 1.0 <- współczynnik eksploracji

Liczba epizodów w każdej symulacji = 100000

Eksperymenty

Badanie wpływu parametru H – horyzontu

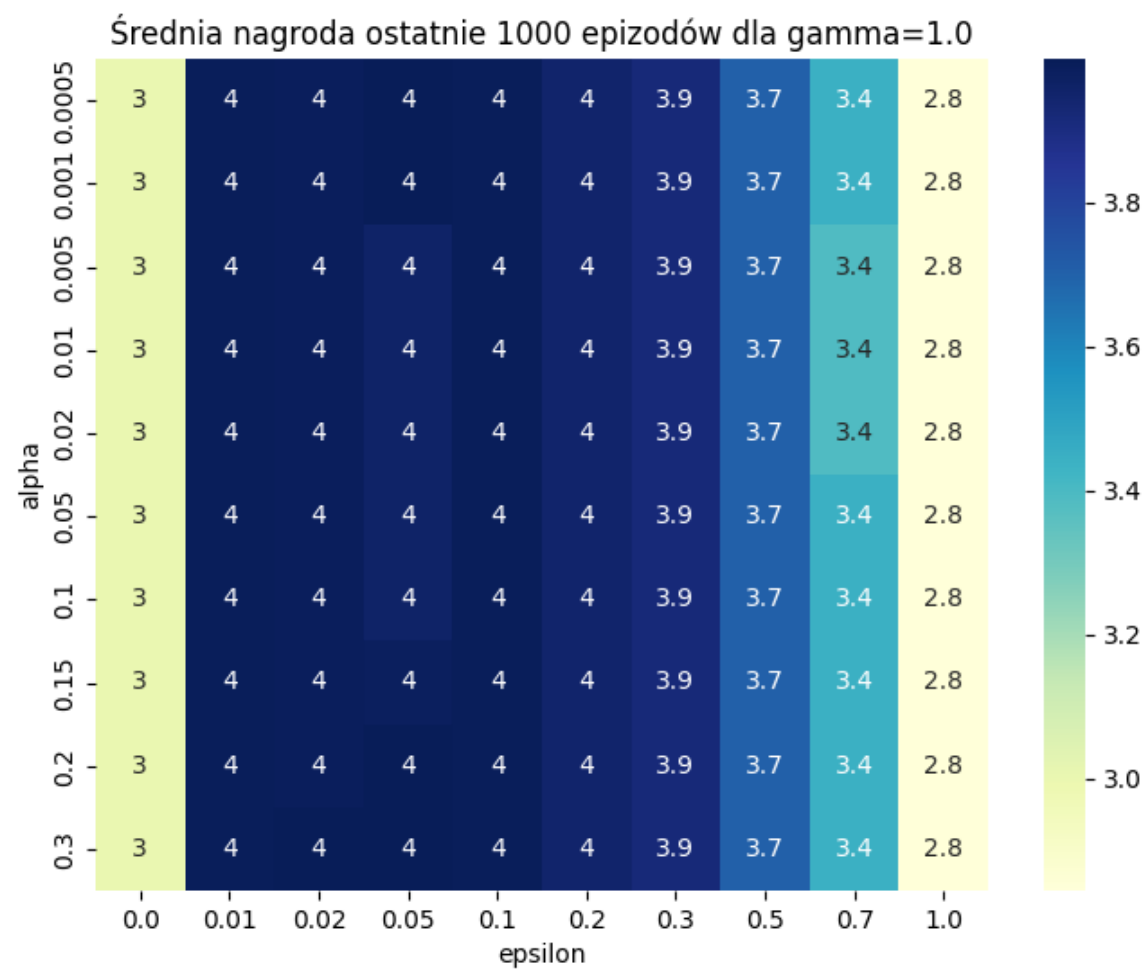
H = 4



a

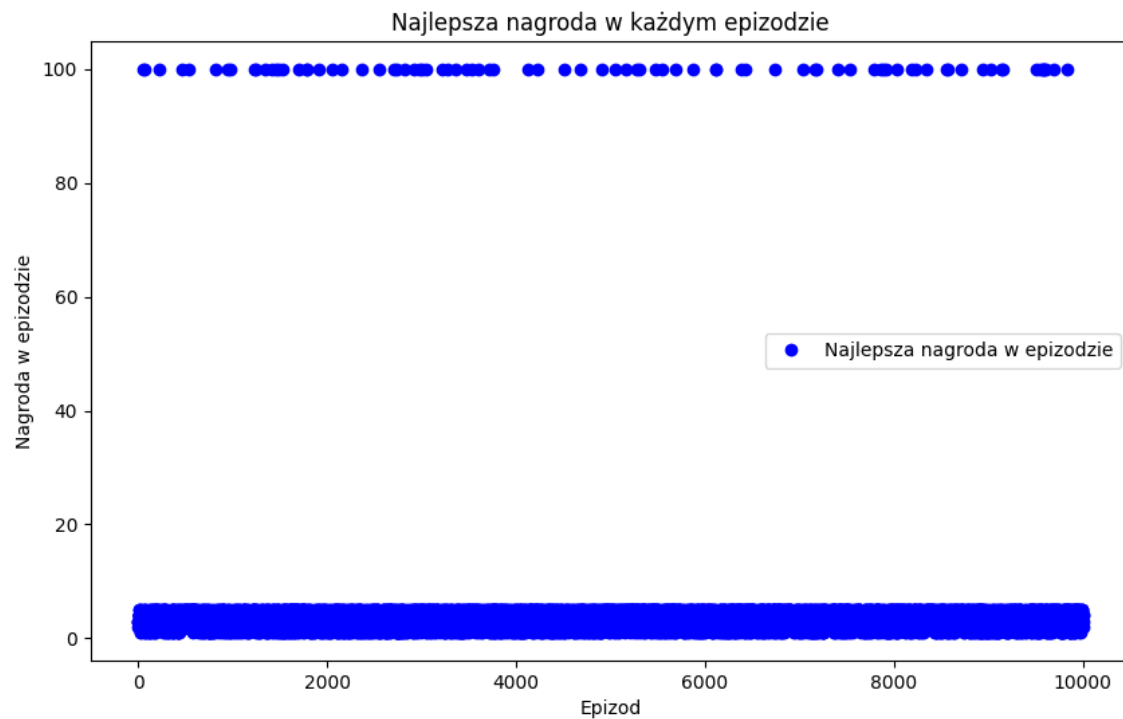
Na powyższej heatmapie, widzimy jakie wartości osiąga agent dla $H = 4$, widać tendencję przy wartościach epsilon bliskich 0, że agent wybiera szybki zysk, zamiast eksploracji i dzięki temu przy małym horyzoncie osiąga najwyższe w tym przypadku nagrody. Z drugiej strony widzimy, że agenci z wyższą tendencją do eksploracji, osiągają średnio gorsze wyniki.

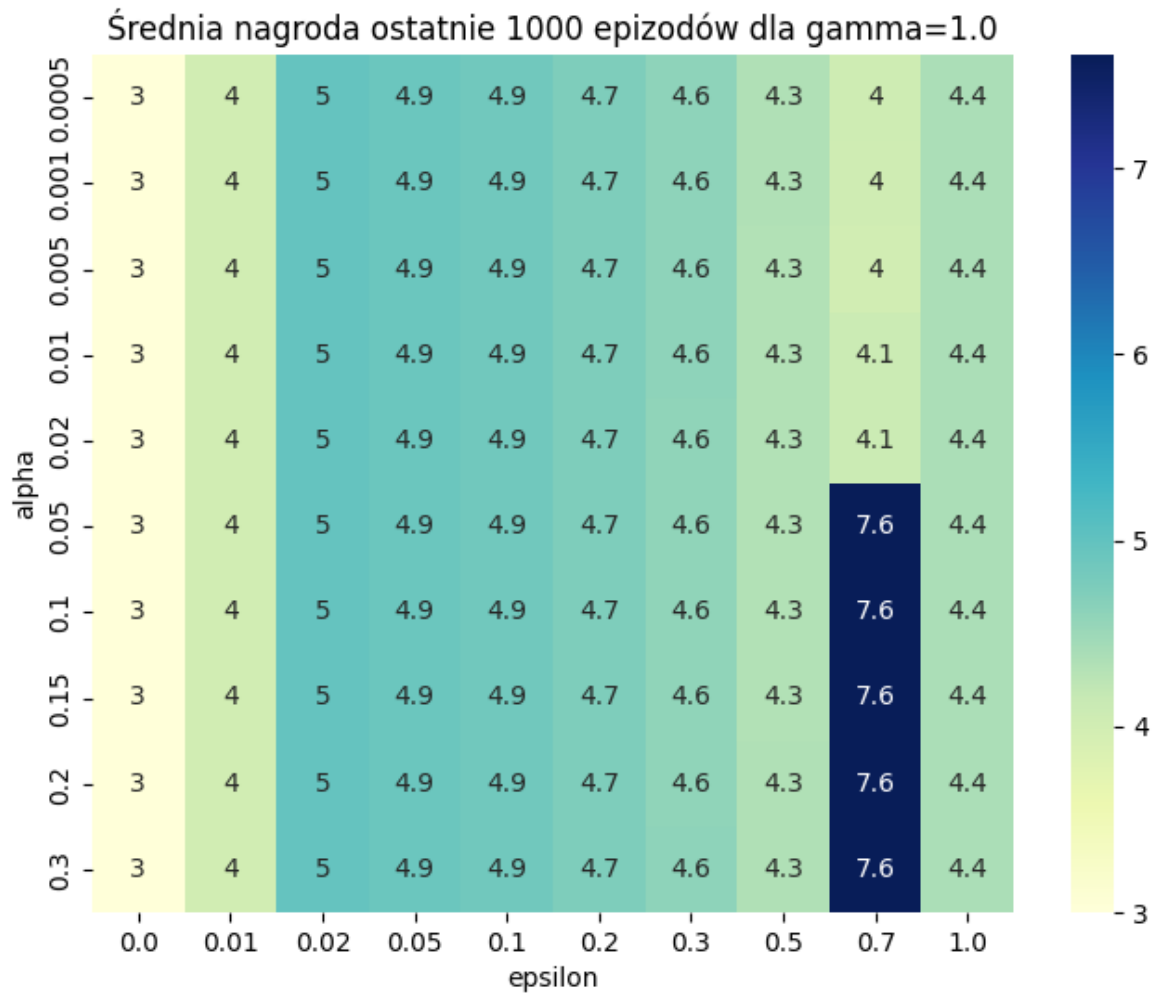
H = 6



W przypadku $H = 6$, widzimy, że agenci nie są w stanie osiągać stanu 10, przez co nadal najlepsze nagrody, daje strategia względnie negatywnie nastawiona do eksploracji. Widzimy, że im agent chętnie eksploruje, tym niższe wyniki osiąga.

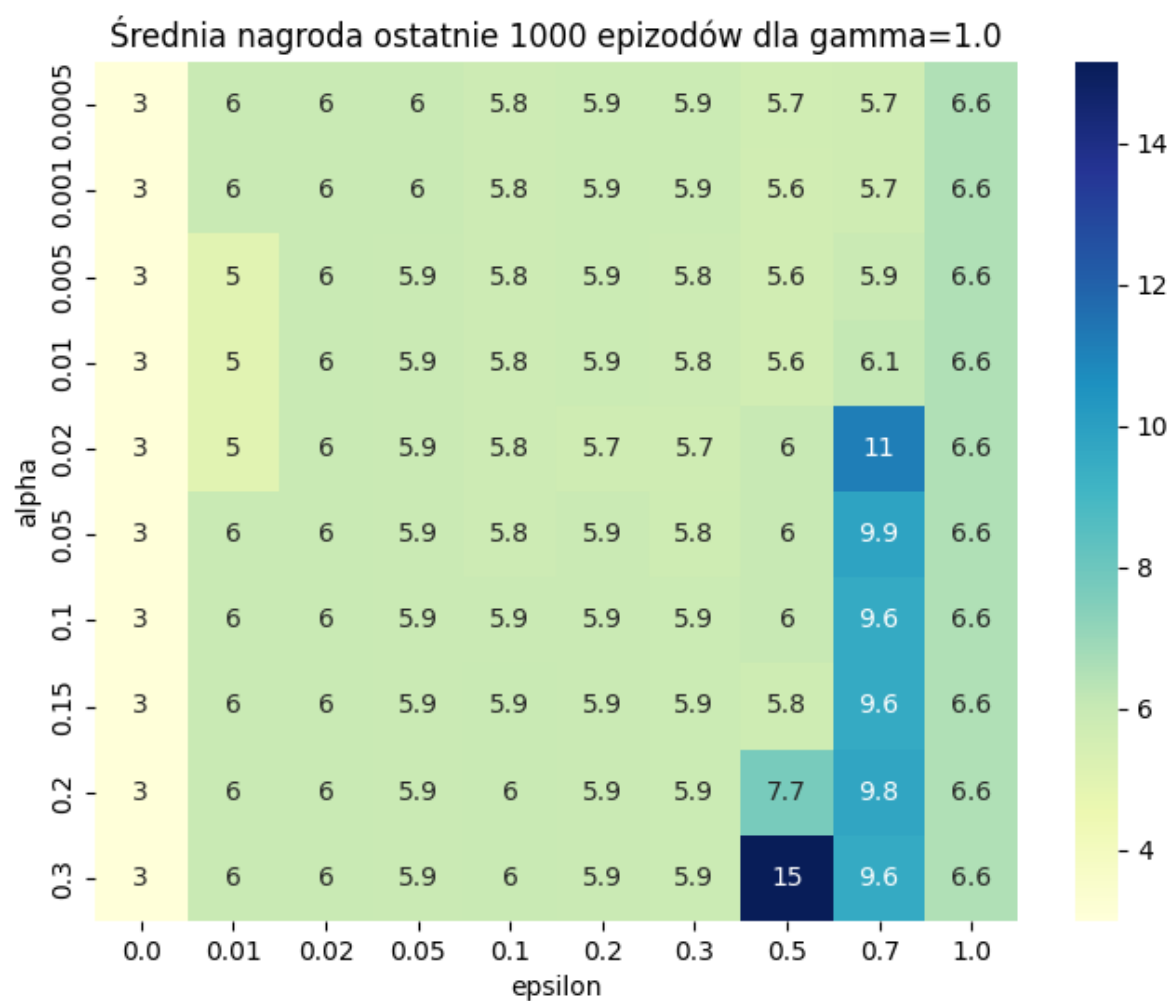
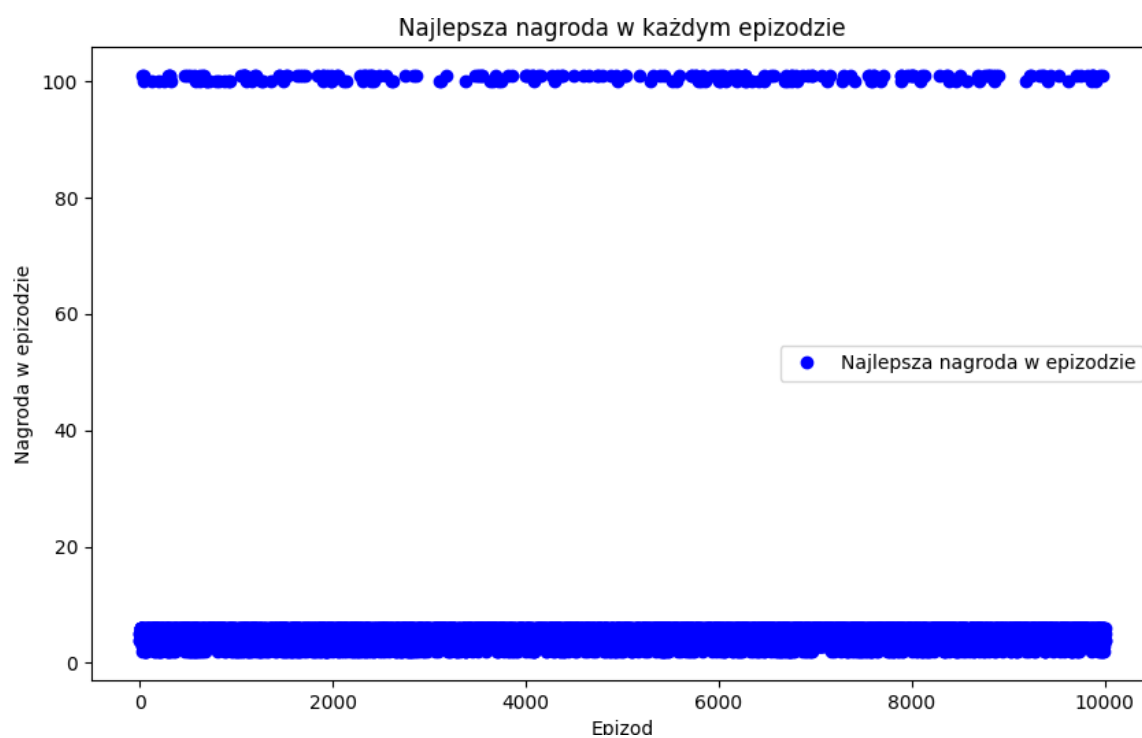
$H = 7$





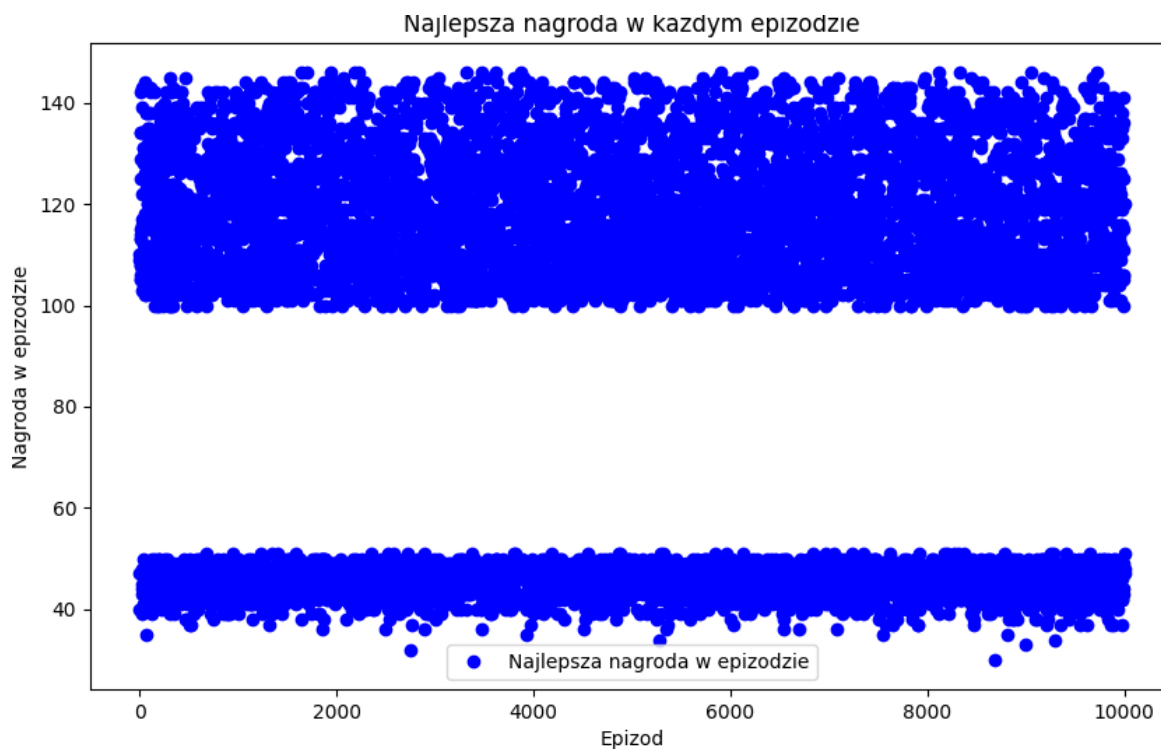
Dla $H = 7$ widzimy, że jest to horyzont przelomowy. Pojedynczy agenci są w stanie osiągnąć wartości nagród > 100 , jednak zdecydowana większość ma nagrody w zakresie 4-5. Na heatmapie, widzimy, że najlepsze wyniki osiągane są dla $\epsilon = 0.7$ i $\alpha 0.3-0.05$. Jednak przy horyzoncie $= 7$, nagroda 7 oznacza 7 akcji sprzedaży co jest niemożliwe, bez żadnych akcji kupno, przez co wiemy, że średnia jest zawyżona przez agentów, którzy osiągnęli nagrodę za wejście do stanu 10.

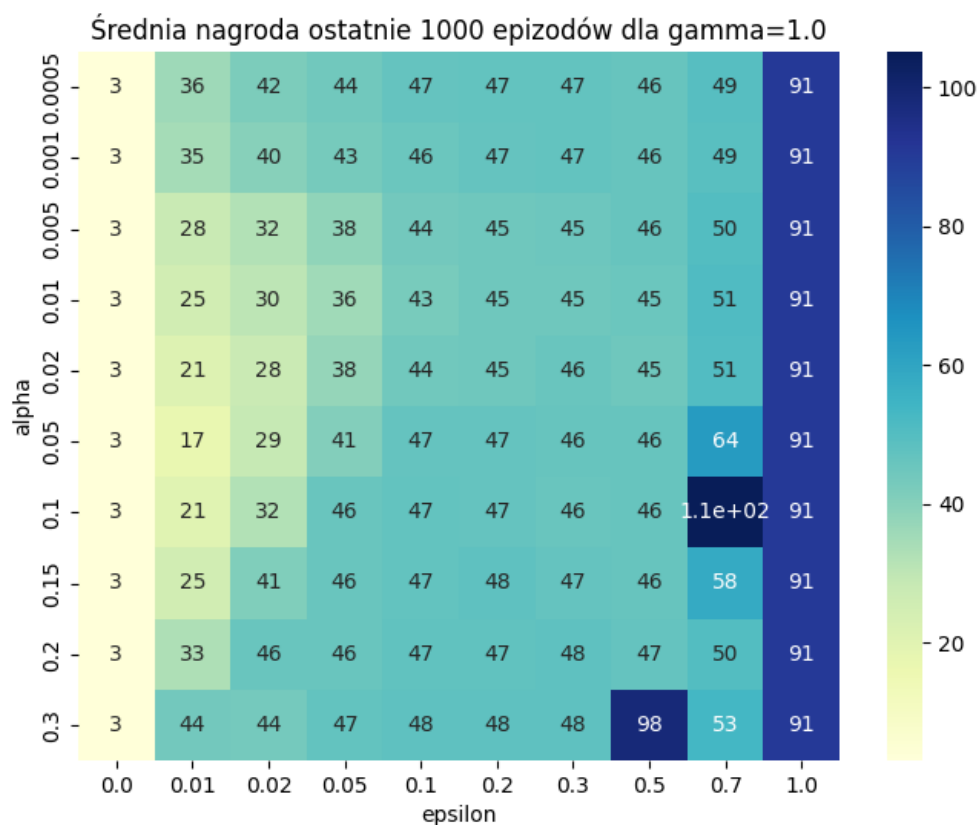
H = 10



W przypadku horyzontu = 10, widzimy większy przekrój wartości nagród osiąganych przez agenta. Teraz tendencja związana z parametrem jest odwrotna, agent stawiający na szybką nagrodę osiąga niskie wyniki. Za to najlepsze wyniki są dla parametru $\epsilon = 0.5$ i $\alpha = 0.3$. Co wskazuje na całkiem wysoką tendencję do eksploracji. Jednak mimo większego horyzontu, średnia jest niska, bo osiągnięcie stanu 10 związanego z nagrodą wymaga 7 akcji kupna. Jednak widzimy, że pojedynczy agenci osiągają rezultaty > 100 , jednak jest to zbyt mało by wpłynąć znacząco na średnią.

H = 100





W przypadku dużego horyzontu np. 100. Widzimy znaczny wzrost średnich wartości nagród osiągniętych przez agentów. Większość agentów osiąga wartości > 100 , z czego największe wartości to około 140-150. Widzimy, też że średnio najbardziej optymalnym parametrem, jest $\epsilon = 1.0$, pozwalający osiągać największe nagrody. Obserwujemy też znikomy wpływ parametru α .

Wnioski

Horyzont ma duży wpływ na to jakie parametry pozwalają agentom osiągać najlepsze wyniki. Mały horyzont będzie premiował agentów z małym parametrem ϵ , przez co agenci będą stawiać na krótkoterminowe zyski. Z drugiej strony, duży horyzont premiuje agentów z dużym ϵ , którzy stawiają na eksplorację, dzięki czemu wiedzą o nagrodzie za $s = 10$.

Przez cały eksperyment widzimy, znikomy wpływ parametru α , mówiącym o learning rate'cie algorytmu. Dla badanych wartości osiągnięte wyniki są zbliżone, z niewielkimi odchyleniami.

Możemy zatem stwierdzić, że w przypadku tego problemu, w celu znalezienia optymalnej strategii dla agenta zarządzającego stanem sklepu, warto dla mniejszych wartości horyzontu ustawić mniejszy ϵ , lecz w celu osiągnięcia większego zadowolenia, parametr ϵ powinien mieć wartość bliską 1.0.