

# Uczenie ze wzmocnieniem

Mateusz Plinta

## Wstęp

Wykonano po 5 symulacji dla każdego z 9-ciu setów parametrów definiujących przebieg symulacji. Krótki opis parametrów:

- **alpha** - stosunek pomiędzy dotychczasową wiedzą a wnioskami z nowego doświadczenia
- **gamma** - współczynnik określający stosunek nagrody już otrzymanej oraz przewidywanych nagród w kolejnych krokach
- **exploration\_rate** - Szansa na wykonanie losowego wyboru(eksperymentu)
- **exploration\_min** - minimalna wartość **exploration\_rate**
- **exploration\_decay** - współczynnik degradacji **exploration\_rate** po każdej iteracji
- **exploration\_degrade** - *<boolean>* definiuje czy degradować **exploration\_rate** czy nie

Przy nauczaniu w słowniku *knowledge* wykorzystywano 0 *reward* gdy nie był to koniec iteracji, oraz ujemny -1, gdy klocek spadał. Dzięki temu uzyskano system który w miarę możliwości szybko się uczy.

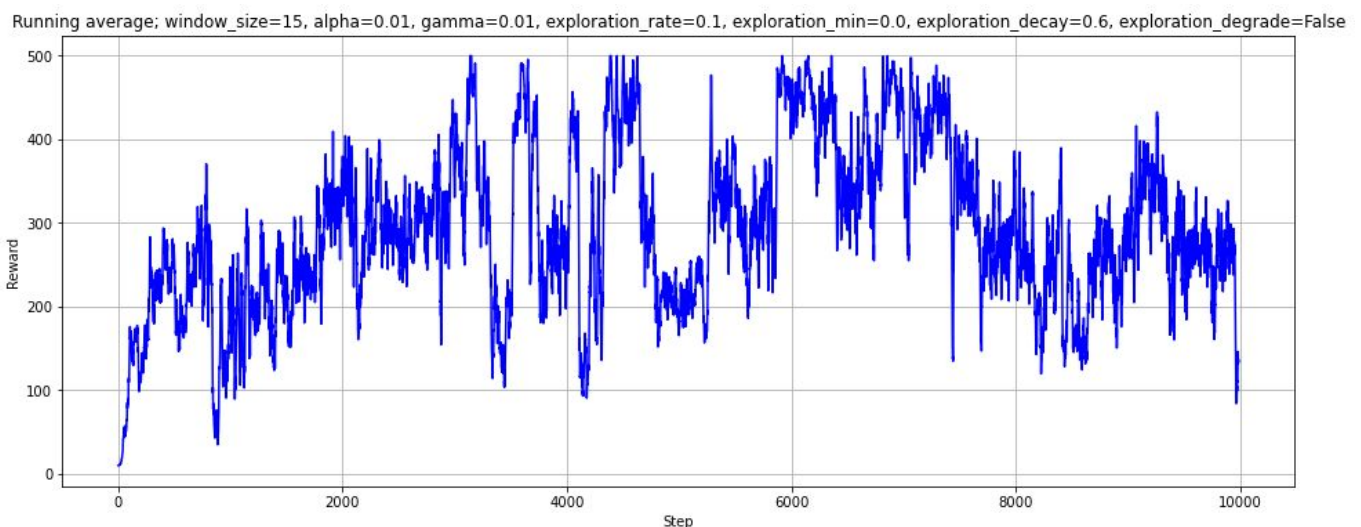
## Zad. 1

- kod w załączniku

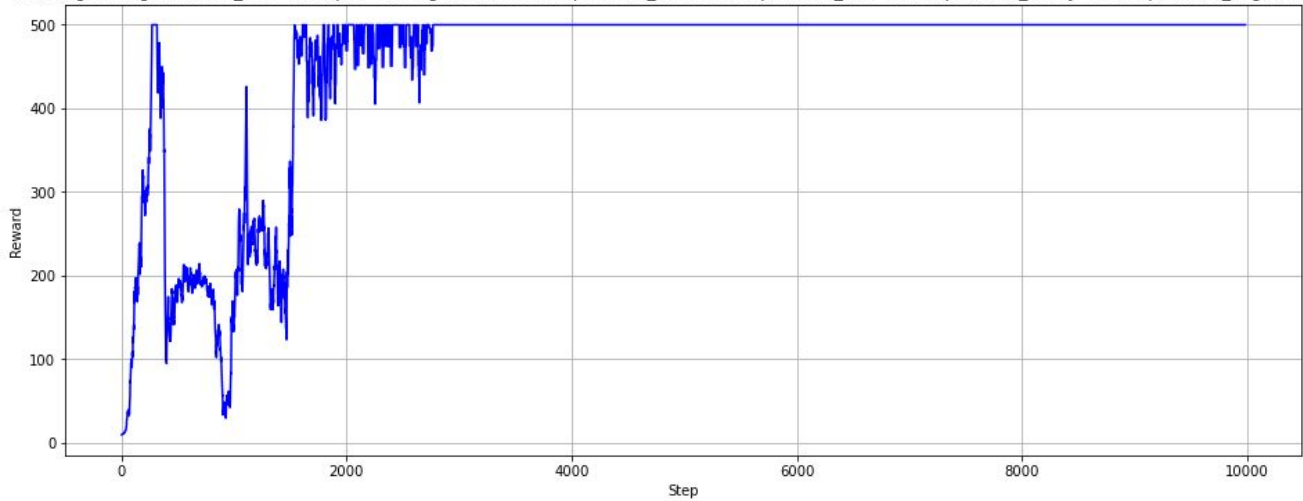
## Zad. 2

Każda symulacja trwała 10000 epok. W przypadku poniższych wykresów przedstawiających średnią kroczącą, wybrano okno szerokości 15 iteracji do uśredniania wyników. W tytule każdego wykresu znajdują się wszystkie parametry symulacji.

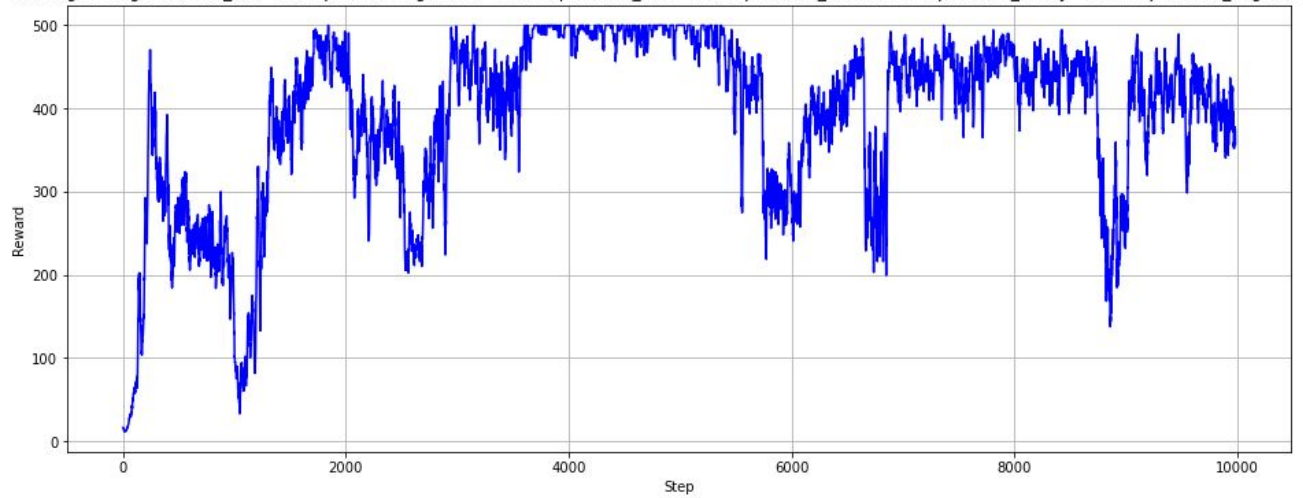
## Średnia krocząca



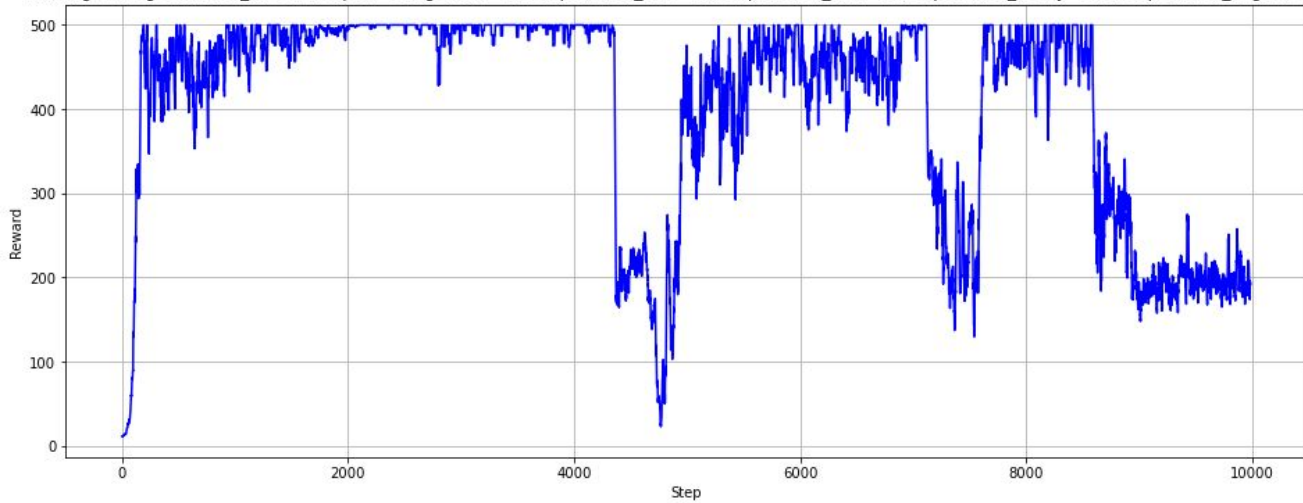
Running average; window\_size=15, alpha=0.01, gamma=0.01, exploration\_rate=0.1, exploration\_min=0.0, exploration\_decay=0.6, exploration\_degrade=True



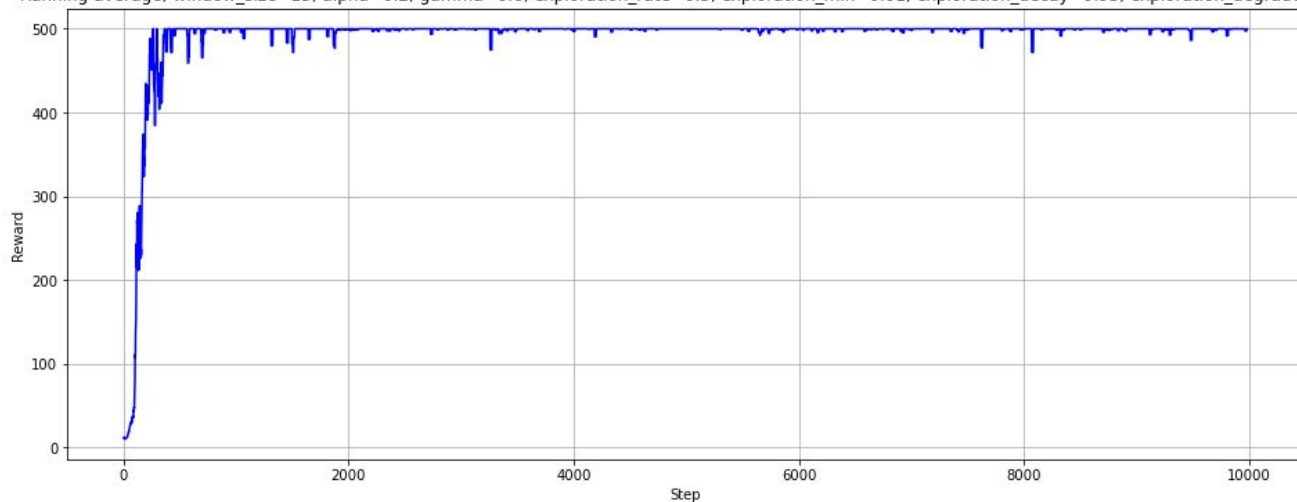
Running average; window\_size=15, alpha=0.01, gamma=0.01, exploration\_rate=0.8, exploration\_min=0.01, exploration\_decay=0.95, exploration\_degrade=True



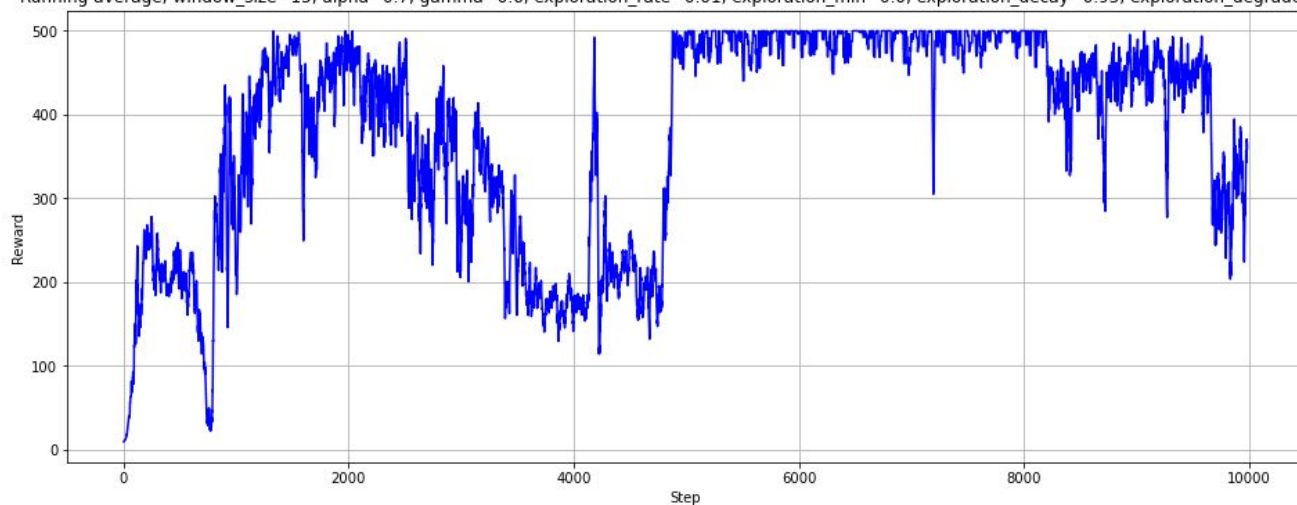
Running average; window\_size=15, alpha=0.1, gamma=0.6, exploration\_rate=0.3, exploration\_min=0.0, exploration\_decay=0.95, exploration\_degrade=True



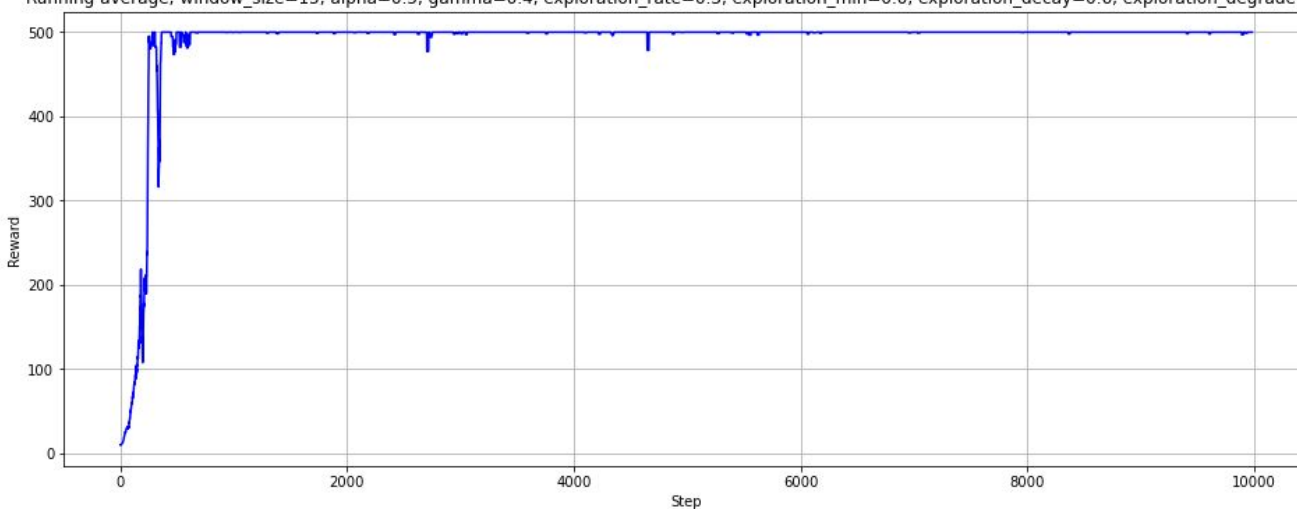
Running average; window\_size=15, alpha=0.2, gamma=0.6, exploration\_rate=0.5, exploration\_min=0.01, exploration\_decay=0.95, exploration\_degrade=True



Running average; window\_size=15, alpha=0.7, gamma=0.6, exploration\_rate=0.01, exploration\_min=0.0, exploration\_decay=0.95, exploration\_degrade=True

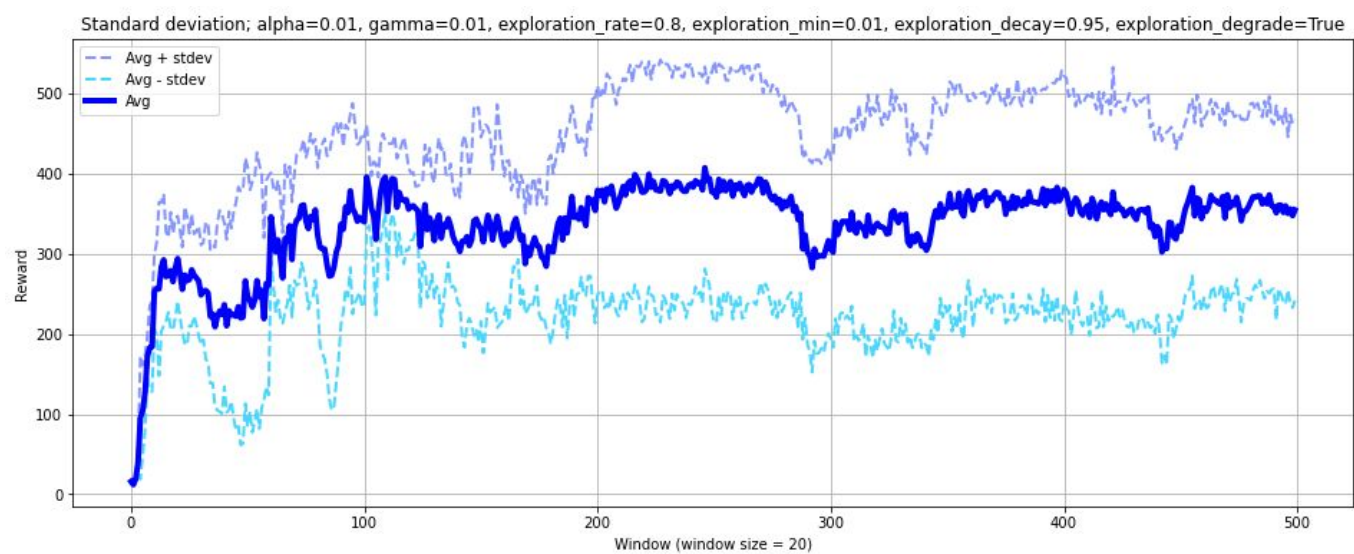
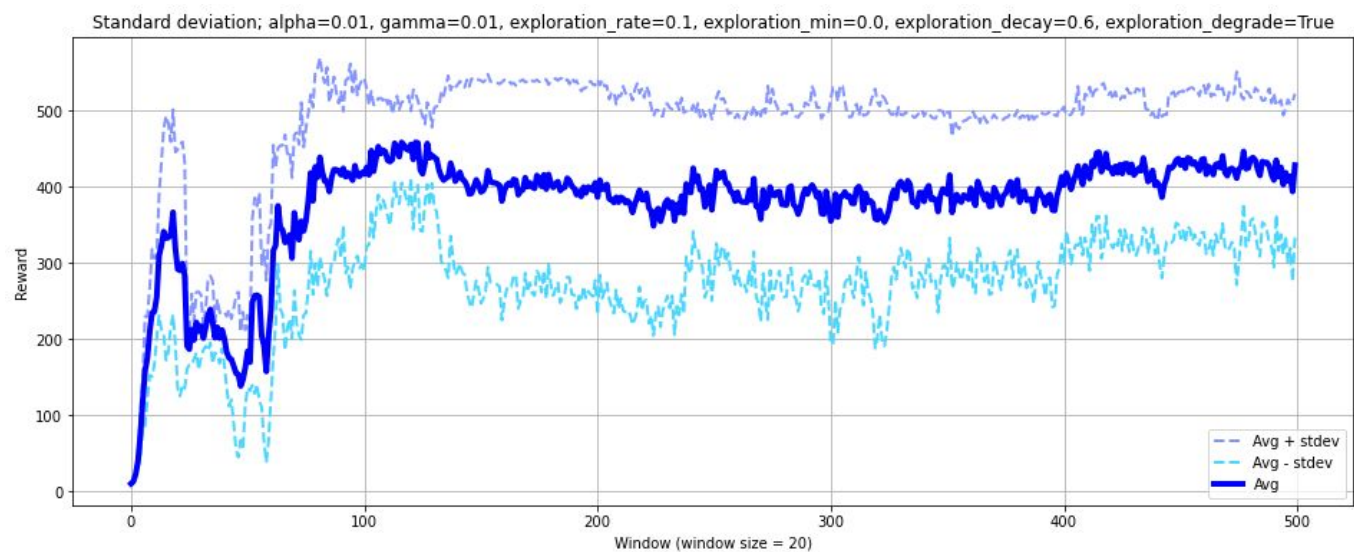
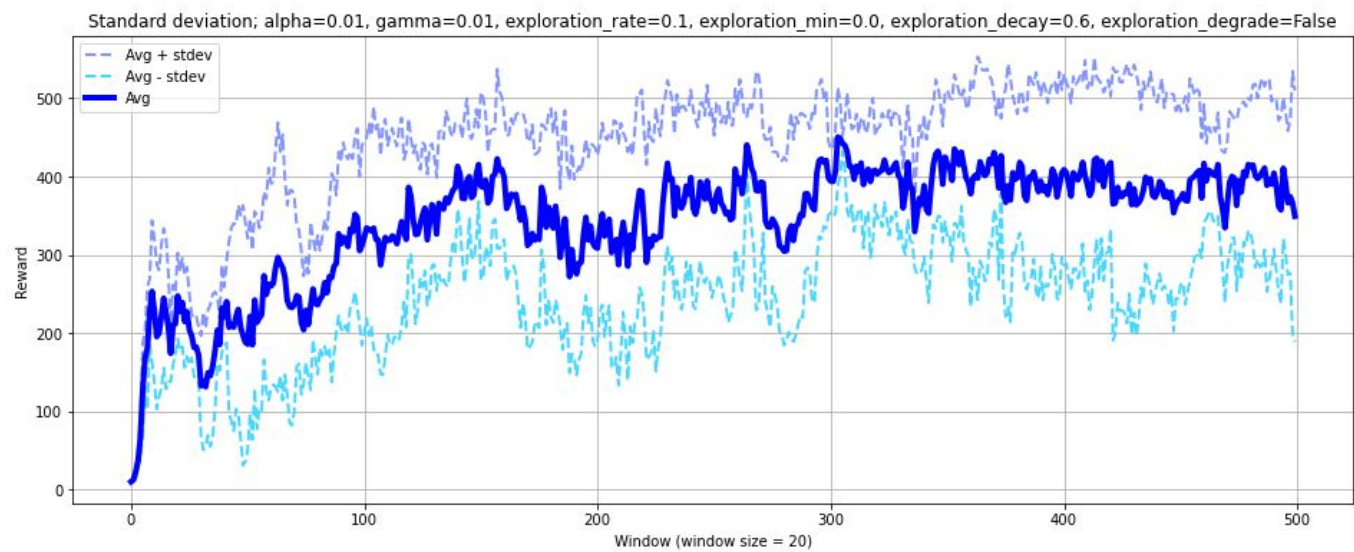


Running average; window\_size=15, alpha=0.5, gamma=0.4, exploration\_rate=0.3, exploration\_min=0.0, exploration\_decay=0.6, exploration\_degrade=True

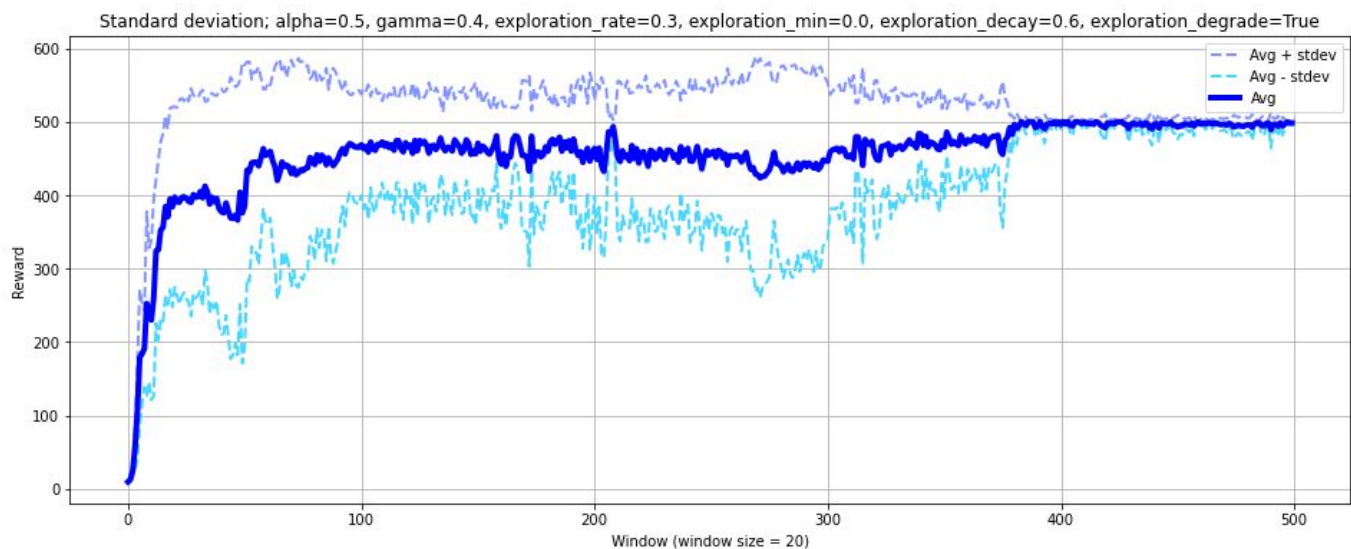
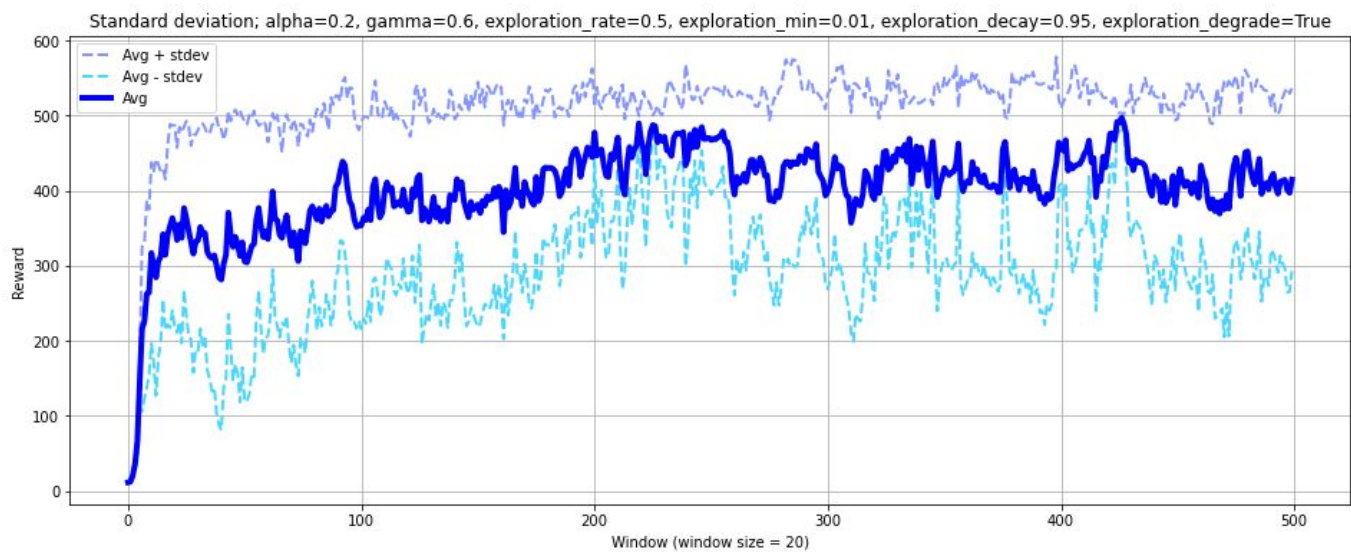
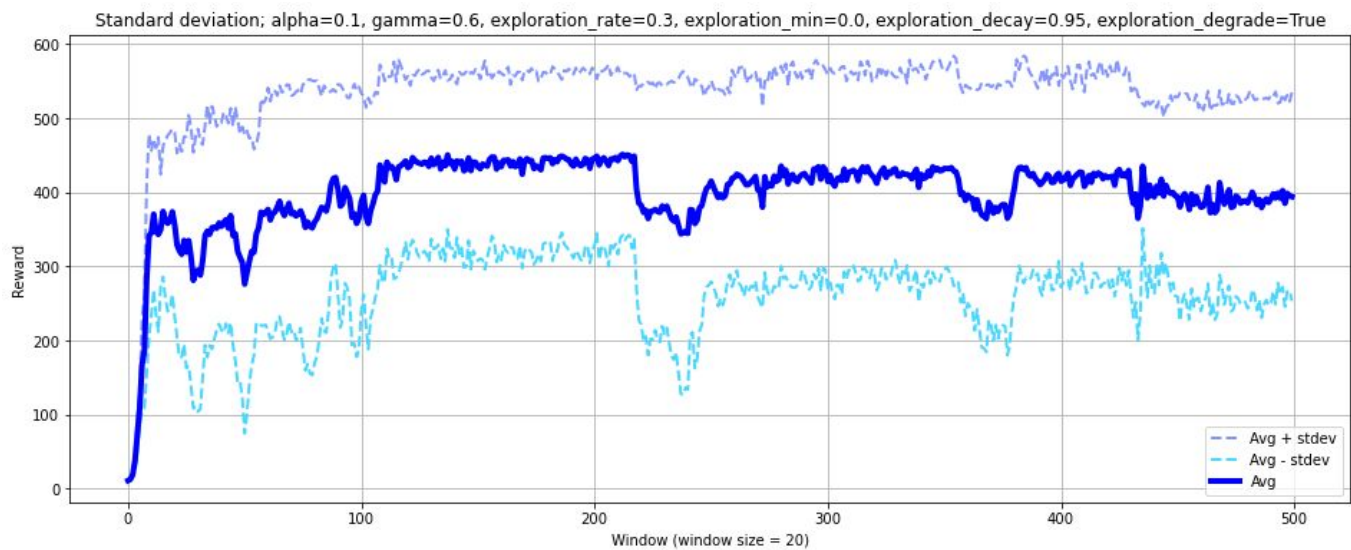


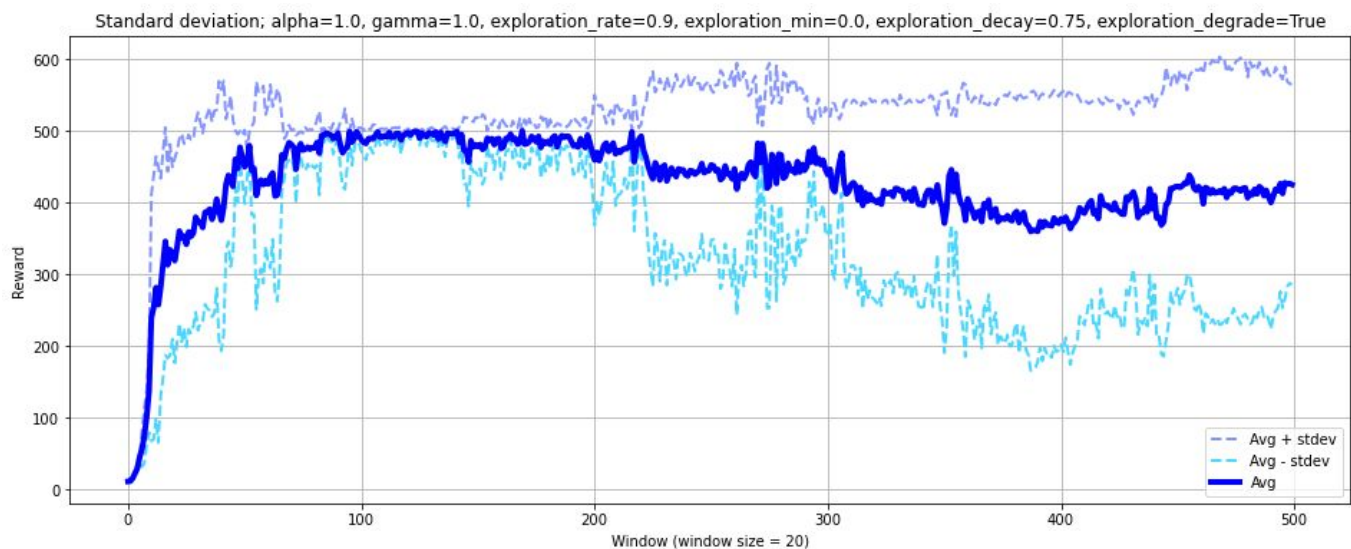
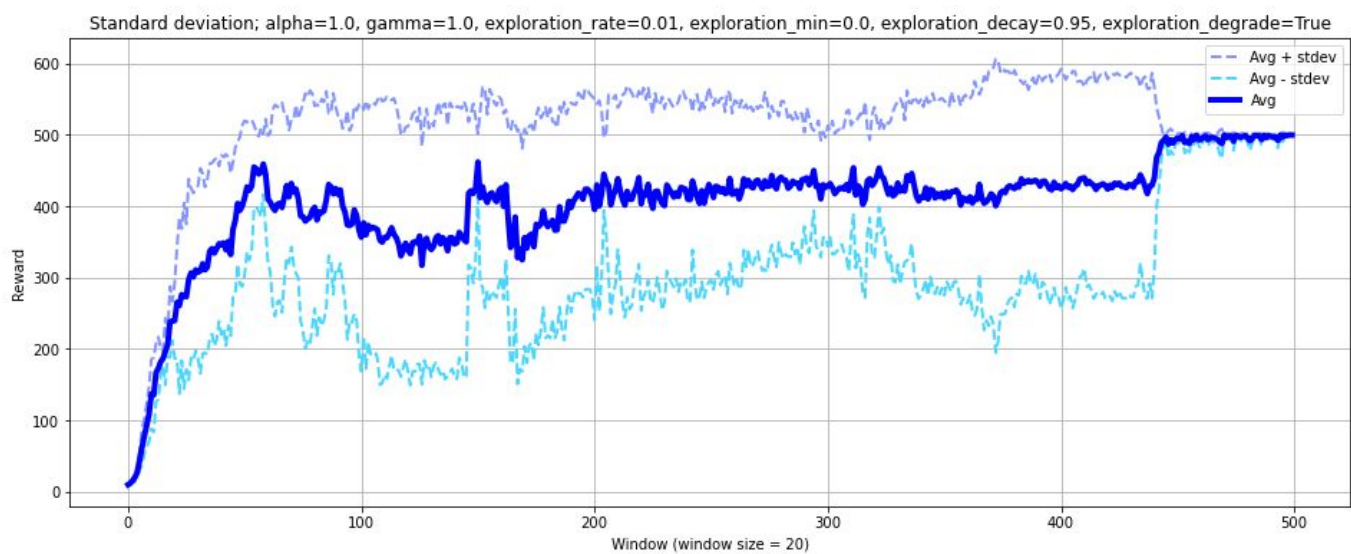
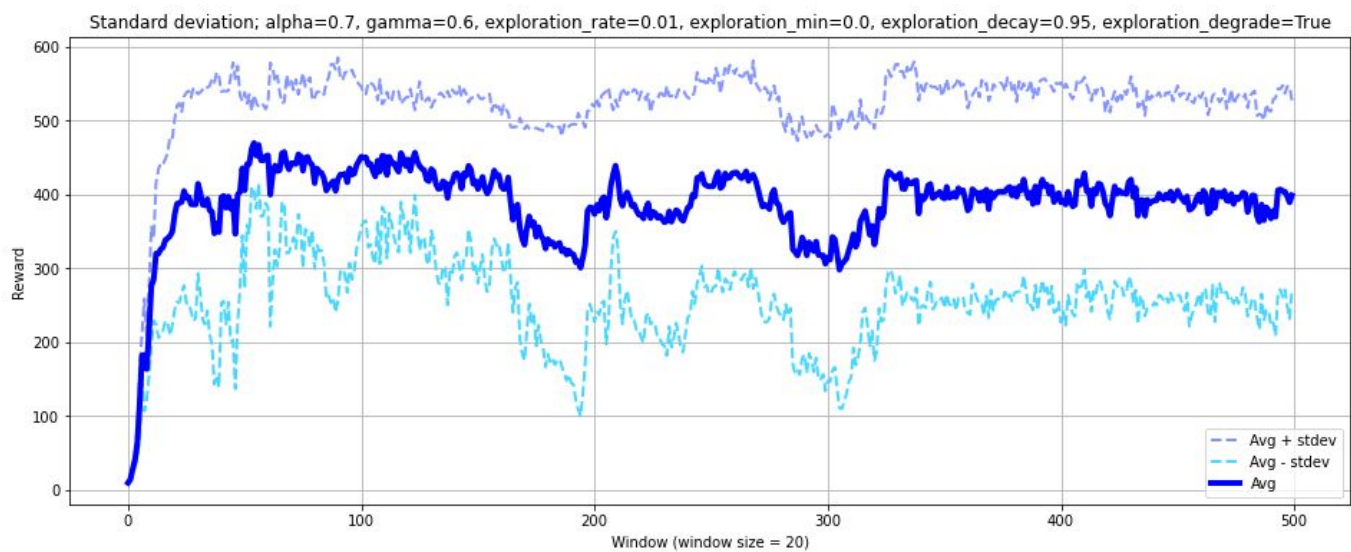
Powyższe wykresy są wyznaczone na podstawie pierwszych ze zbioru pięciu prób symulacji o tych samych parametrach.

# Odchylenie standardowe









Powyższe wykresy uśredniono zakładając okno wielkości 20 iteracji. Linie przerywane przedstawiają wartości wykresu +/- odchylenie standardowe wartości w danym oknie.

## Wnioski

Pierwszy wykres przedstawia symulacje z wyłączoną degradacją **exploration\_rate**. Wykres poniżej zakłada te same parametry, lecz z włączonym tym parametrem. Widać wyraźnie zwiększone zaszumienie pomiarów przy wyłączonym **exploration\_degrade**, spowodowane nie zmniejszaniem się wartości **exploration\_rate**. Mimo tego jednak wygląda na to że z wyłączonym parametrem **exploration\_degrade** system uczy się wolniej, lecz stabilniej. Ostatecznie jednak oba wykresy kończą na średnio podobnych wynikach ~400 score.

Najlepsze wyniki uzyskano dla parametrów:

```
alpha=0.5, gamma=0.4, exploration_rate=0.3, exploration_min=0.0,  
exploration_decay=0.6, exploration_degrade=True
```

gdzie dane rosną bardzo szybko do ~450 score, w miarę stabilnie. Po ok. 400 oknach system stabilizuje się na 500. Mogły na to mieć wpływ stosunkowo duże wartości parametrów **alpha** i **gamma** oraz **exploration\_rate**, który stosunkowo duży na początku pozwala na eksplorację potencjalnie lepszych możliwości.

Dosyć dobre rezultaty uzyskano również wykorzystując poniższy zestaw parametrów:

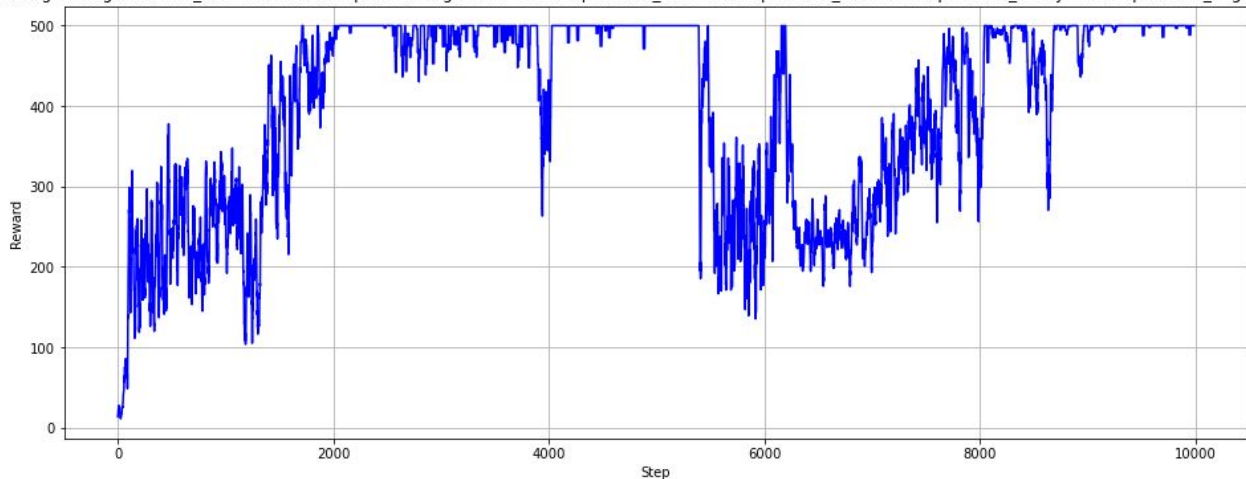
```
alpha=1.0, gamma=1.0, exploration_rate=0.01, exploration_min=0.0,  
exploration_decay=0.95, exploration_degrade=True
```

gdzie score również szybko rośnie do wartości 400, jednakże zdecydowanie mniej stabilnie niż w poprzednim przypadku. Prawdopodobnie spowodowane jest to przez zbyt duże wartości **alpha** oraz **gamma**. Ostatecznie wynik stabilizuje się na 500 w okolicach 450 okna.

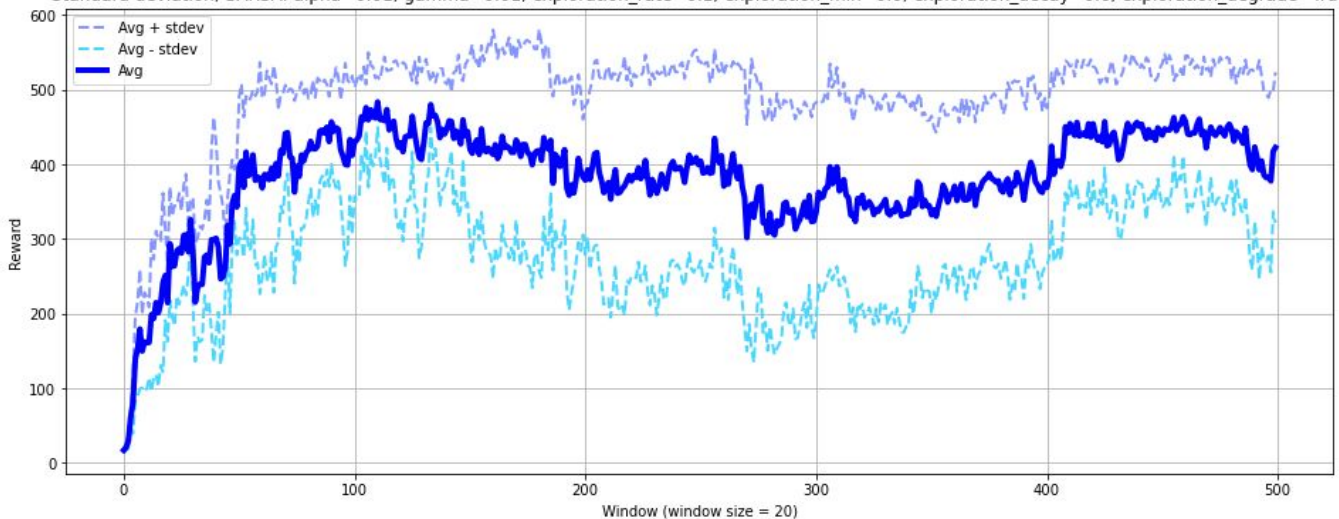
Pozostałe zbiory wartości parametrów powodują końcowo, średnio wyniki rzędu wartości ~400.

## Zad. 3 - SARSA

Running average; window\_size=15, SARSA: alpha=0.01, gamma=0.01, exploration\_rate=0.1, exploration\_min=0.0, exploration\_decay=0.6, exploration\_degrade=True



Standard deviation; SARSA: alpha=0.01, gamma=0.01, exploration\_rate=0.1, exploration\_min=0.0, exploration\_decay=0.6, exploration\_degrade=True



### Poprzednia metoda vs SARSA

Porównując wykresy dla tych samych wartości parametrów dla obu metod, można stwierdzić iż metoda SARSA w moim przypadku powoduje wzrost score'u w podobnym czasie, jednak stan ten się nie utrzymuje - jest o wiele mniej stabilny, przez co w okolicy 200 okna następuje degradacja wyniku. Jest on później odbudowany. W przypadku poprzedniej metody uczenia ze wzmocnieniem obserwujemy większą stabilność wyniku.