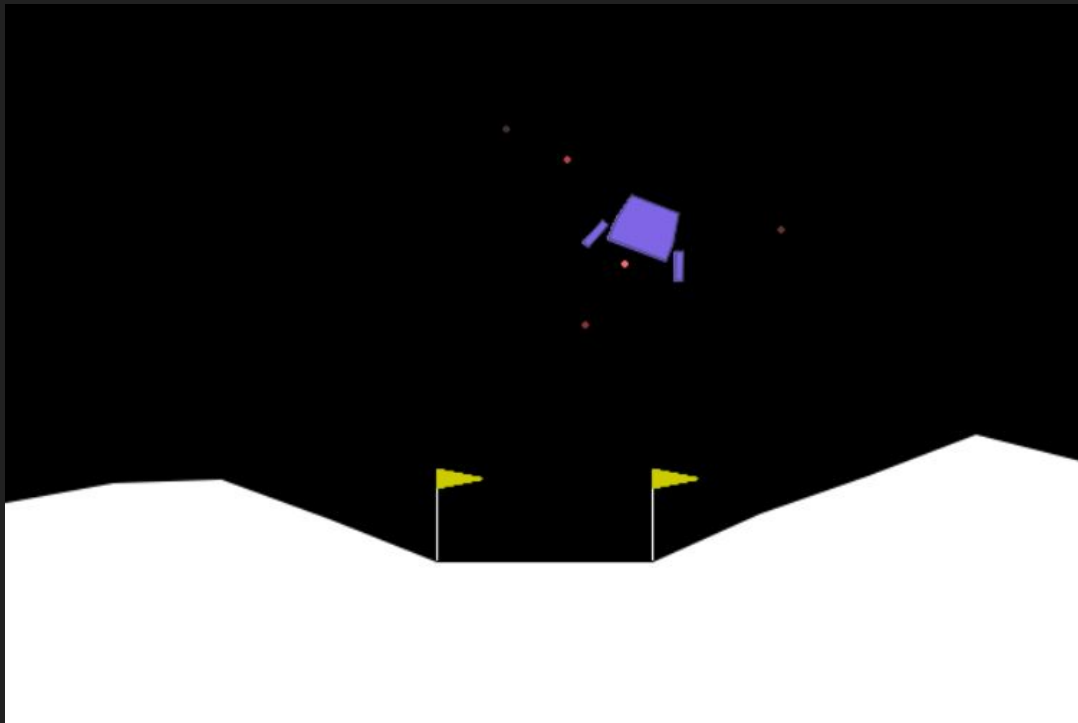


# Q-learning i deep Q-learning na przykładzie Lunar Lander

Krzysztof Surówka

# Lunar Lander



# Środowisko i języki programowania

- Środowisko to Lunar Lander z biblioteki gymnasium (pochodna biblioteki gym od OpenAI).
- Język programowania to Python.
- Do sieci neuronowych została użyta biblioteka PyTorch.
- Dodatkowo do obliczeń numerycznych została użyta biblioteka NumPy.

# Stan w Lunar Lander

Stan składa się z:

- Pozycji  $x$  oraz  $y$ .
- Prędkości  $x$  oraz  $y$ .
- Kąta względem pionu i prędkości kątowej.
- Wartość 0/1 mówiącej czy lewa/prawa noga dotyka ziemi.

Łącznie daje to 8 parametrów.

# Akcje w Lunar Lander

Są 4 możliwe akcje do wykonania w Lunar Lander:

- Uruchomienie głównego silnika.
- Uruchomienie prawego/lewego silnika.
- Brak uruchomienia silnika.

Silnik zostaje uruchomiony na pełnej mocy albo wcale. Istnieje również ciągła wersja tego problemu, gdzie możemy sterować również mocą silnika.

# Nagrody w Lunar Lander

Nagrody w Lunar Lander są liczone w następujący sposób:

- Jest zwiększana/zmniejszana im bliżej/dalej lądownik jest od platformy.
- Jest zwiększana/zmniejszana im wolniej/szybciej porusza się lądownik.
- Jest zmniejszana w miarę większego nachylenia lądownika.
- Jest zwiększana o 10 punktów za każdą nogę, która jest w kontakcie z ziemią.
- Jest zmniejszana o 0,03 punkta za każdą klatkę, w której silnik boczny jest aktywowany.
- Jest zmniejszana o 0,3 punkta za każdą klatkę, w której silnik główny jest aktywowany.

Dodatkowo lądownik otrzymuje nagrodę -100 lub +100 punktów za zderzenie lub bezpieczne lądowanie.

# Q-learning i deep Q-learning

W Q-learningu zostało użyte równanie:

$$Q(s, a) = (1 - \alpha) * Q(s, a) + \alpha * [R(s, a) + \gamma * \max(Q(s', a'))]$$

Natomiast w deep Q-learningu zostało użyte równanie Bellmana:

$$Q(s, a) = R(s, a) + \gamma * \max(Q(s', a'))$$

# Dyskretyzacja stanów

Aby zastosować Q-learning, dziedzinę stanów należy zdyskretyzować.

- Im większa dyskretyzacja, tym model będzie się szybciej uczył, ale może potem okazać się że będzie mało dokładny.
- Dodatkowo pojawia się problem wielkości Q-table. Im na mniejsze przedziały podzielimy dziedzinę, tym model będzie więcej zajmował miejsca.

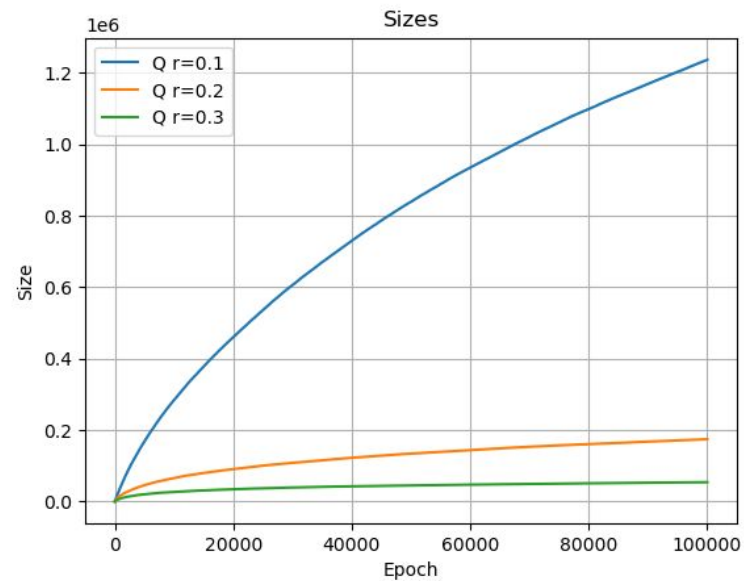
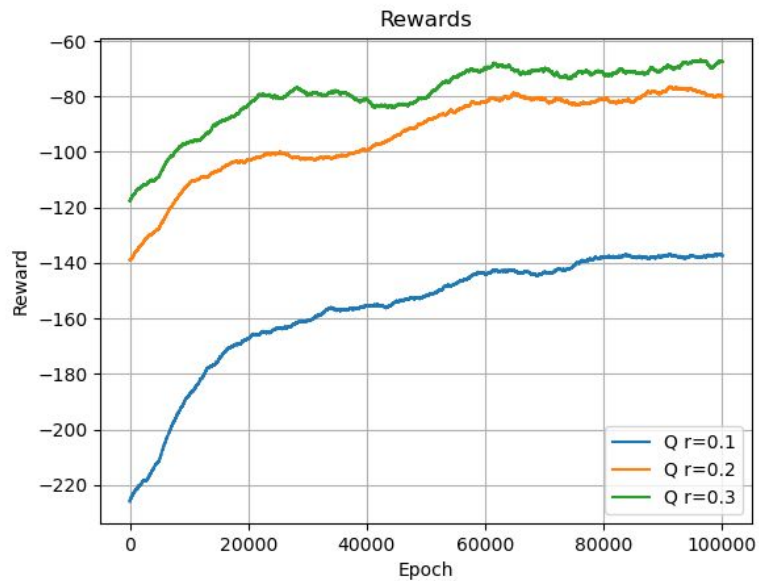
W deep Q-learning nie trzeba dyskretyzować dziedziny, bo nie mam Q-table, a sieć neuronową



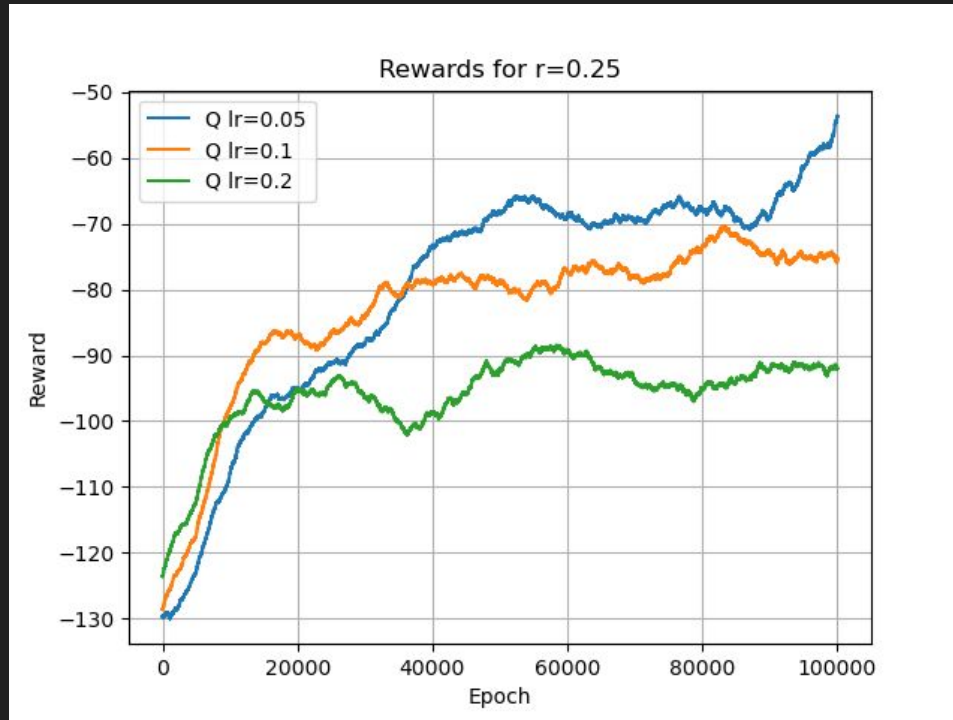
# Inne parametry

- Parametr  $r$  (podział dziedziny na kubelki o rozmiarze  $r$ )
- Parametr  $lr$  (learning rate, parametr  $\alpha$ )
- Parametr  $\gamma$  (discount dla przyszłego stanu,  $\gamma$ )
- Parametr  $u$  (kubelki nie są co  $r$ , a są co  $r^u$ )
- Parametr  $d$  (losowość podejmowania akcji co epokę jest mnożona przed  $d$ , aż do osiągnięcia poziomu 0.05)
- Parametr  $n$  (wielkość warstwy ukrytej w sieci neuronowej)

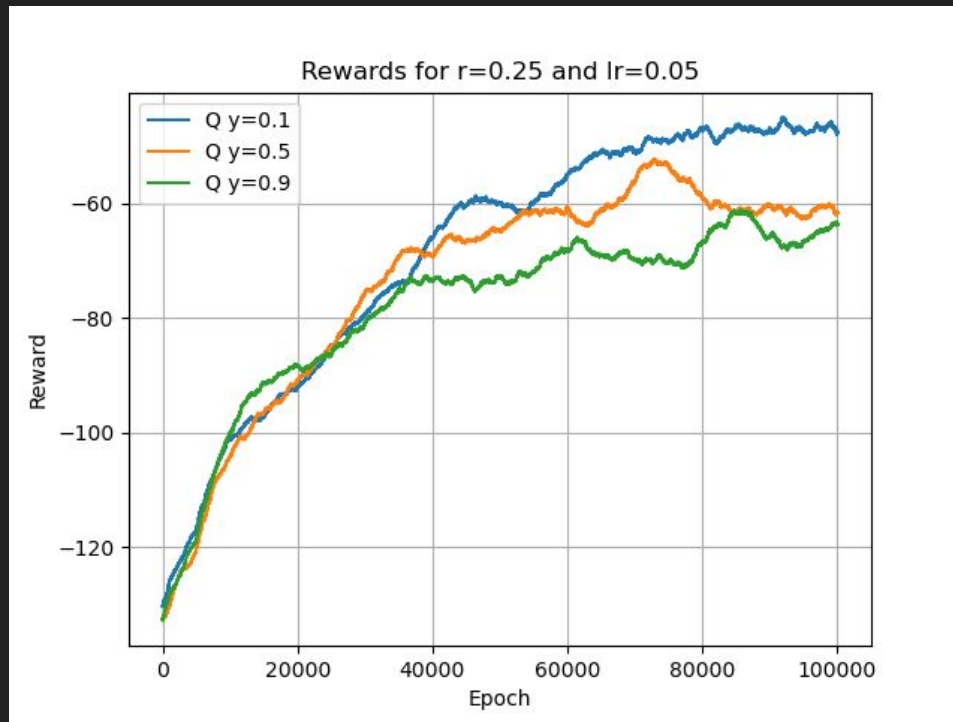
# Parametr r



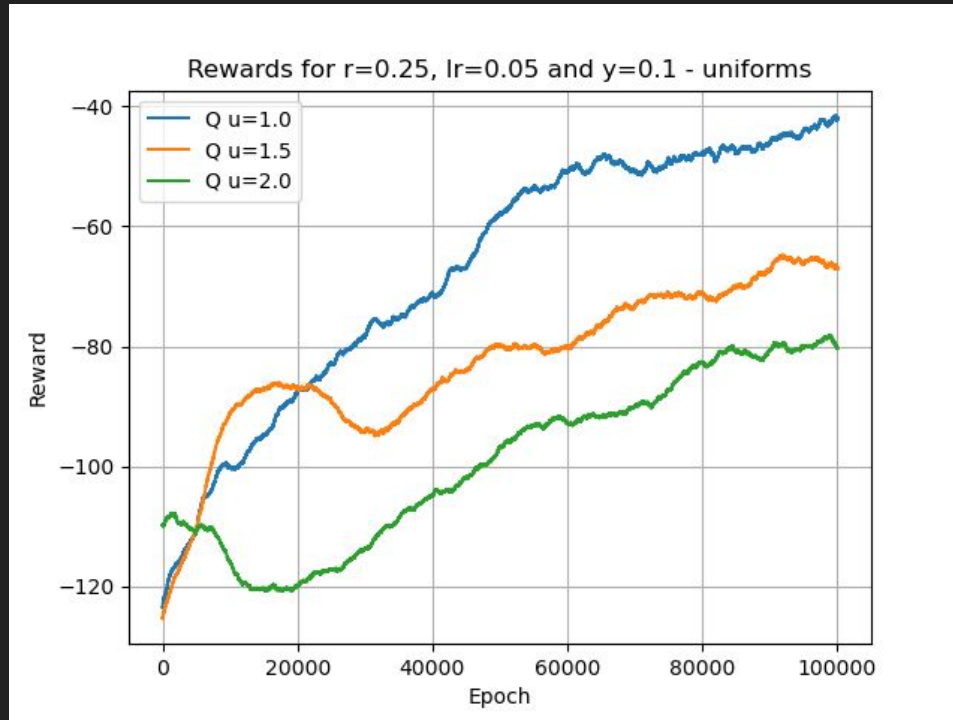
# Parametr lr



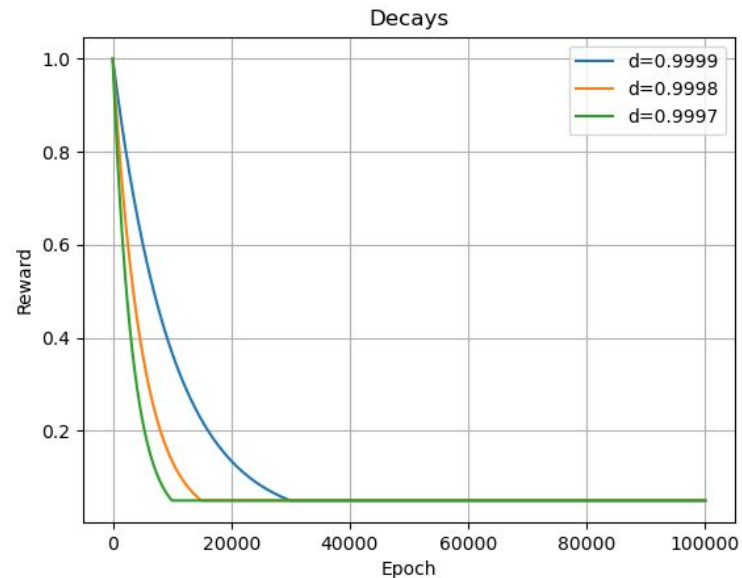
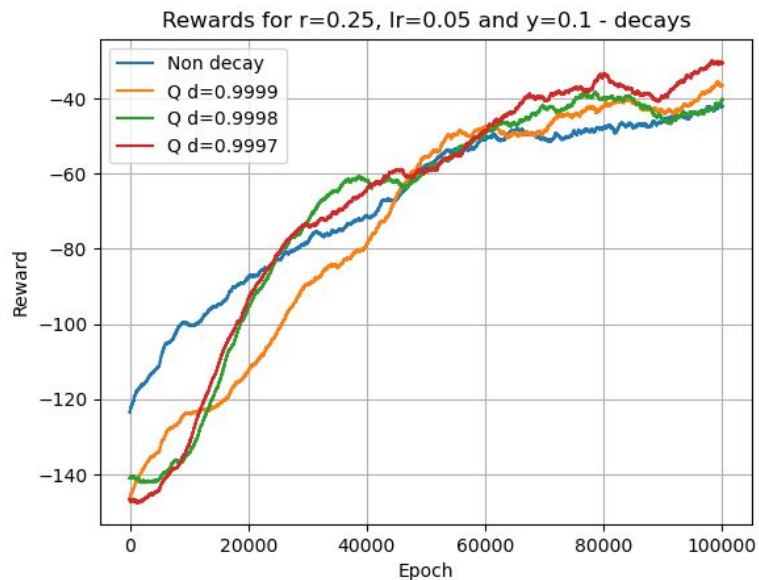
# Parametr $\gamma$



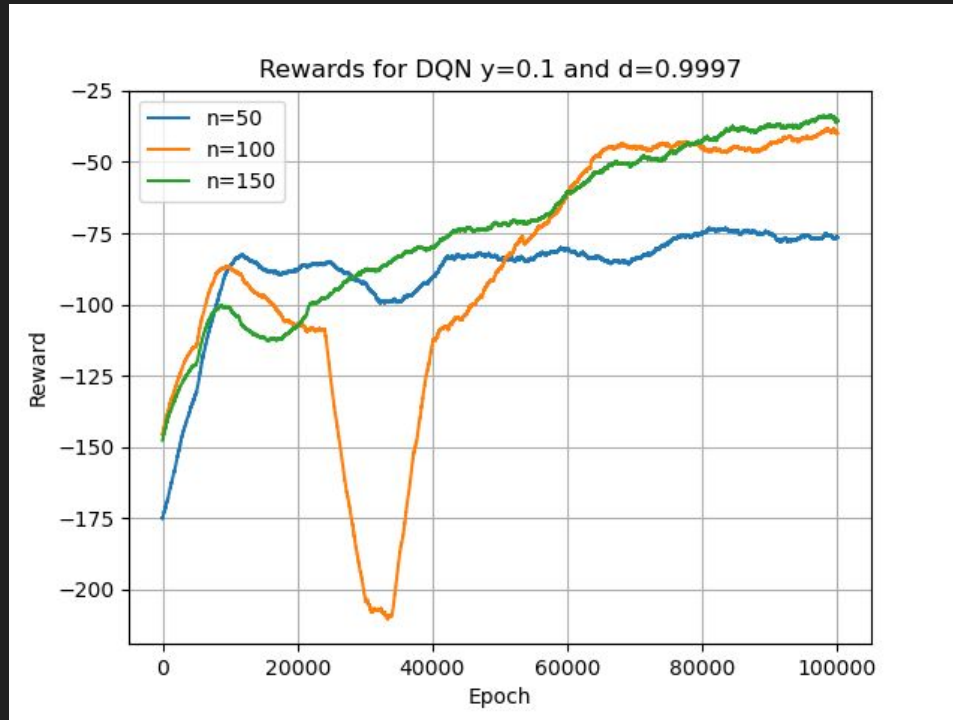
# Parametr $u$



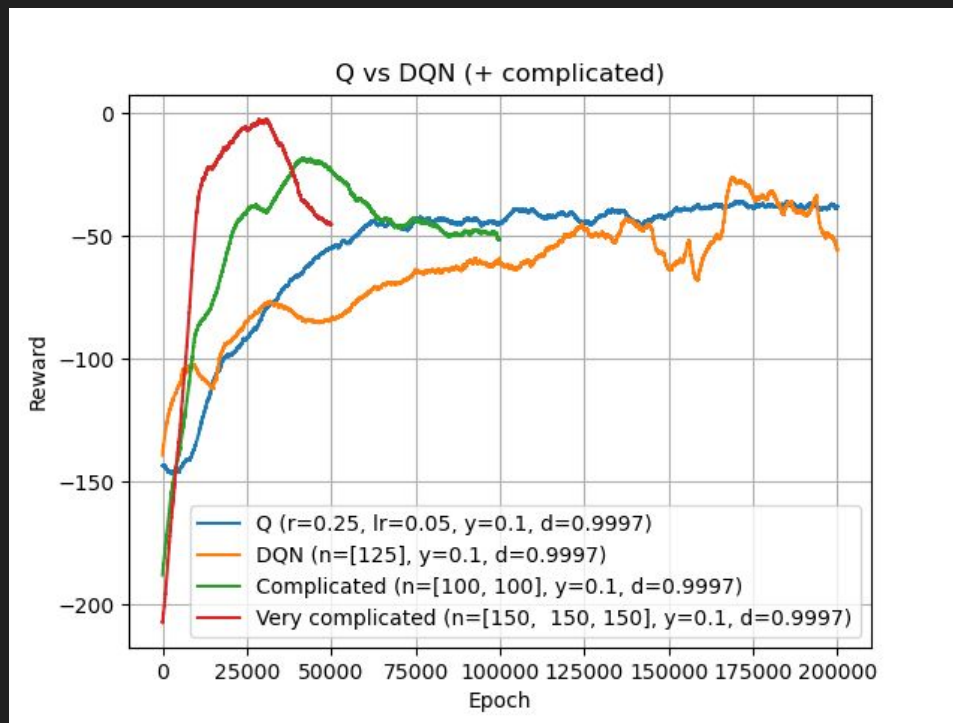
# Parameter d



# Parametr n

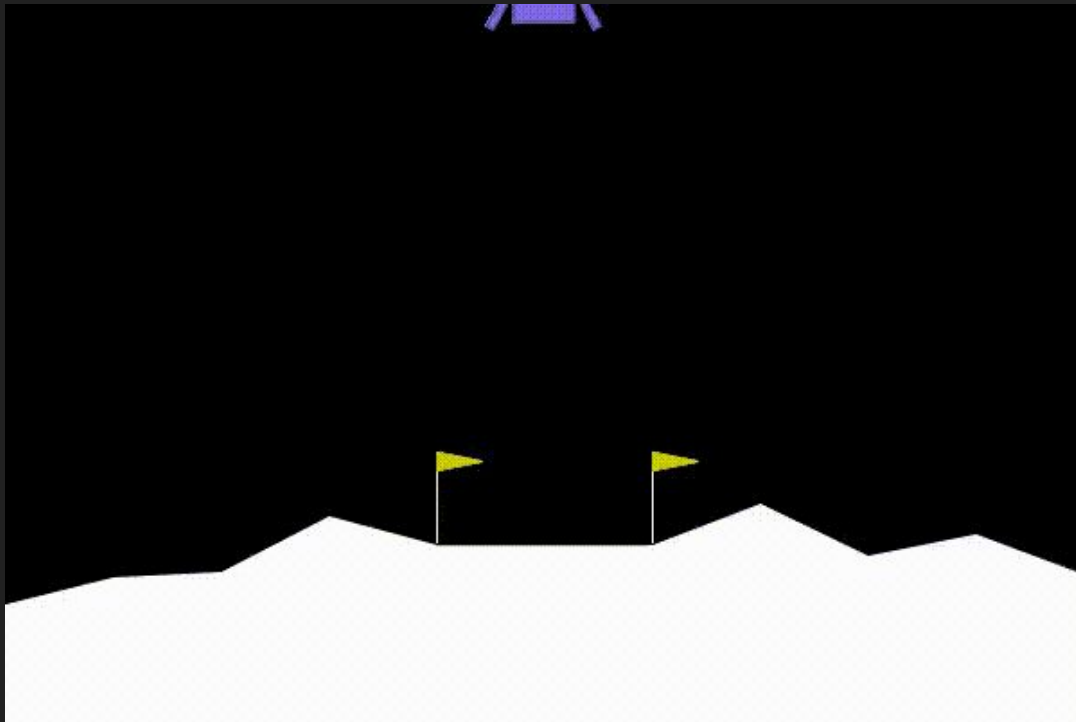


# Q-learning vs deep Q-learning

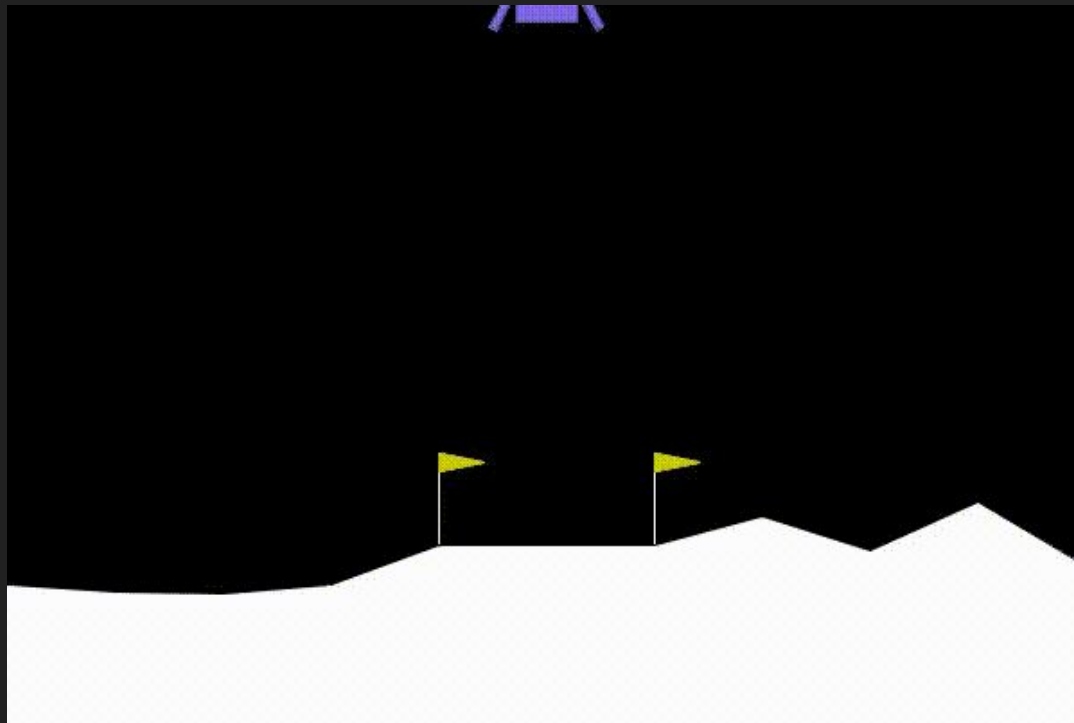




# Rezultat



Rezultat bez kary za paliwo



# Czasy nauk

- Q-learning: 1h / 100000 epok
- Deep Q-learning: 5h / 100000 epok
- Deep Q-learning z dużym modelem: 10h / 100000 epok lub więcej

Dużym problemem w deep Q-learning jest ciągle wywoływanie funkcji forward co każdą klatkę gry, nawet do 500 razy na epokę.

## Przydatne linki:

- [https://gymnasium.farama.org/environments/box2d/lunar\\_lander/](https://gymnasium.farama.org/environments/box2d/lunar_lander/)
- <https://github.com/MrKrisuuu/ReinforcementLearning>