

Projekt 2: Lunar Lander

Dawid Roszman

Czerwiec 2024

1 Opis zadania

Chcemy pomyślnie wyladować sonda na księżycu. Napisz program, który to wykona. Do rozpatrzenia sa trzy możliwości:

- własny skrypt,
- algorytm reinforcement learning (w tym także z siecia neuronowa)
- kontroler rozmyty (wymaga napisania bazy i reguł)
- strategie metaheurystyczne (PSO lub Algorytm Genetyczny)

Na końcu zademonstruj działanie algorytmów na animacjach (możesz zapisać jaki video lub gify)

2 Pierwsze podejście - własny skrypt

W swoim pierwszym podejściu napisałem skrypt, który na podstawie nagrody otrzymanej za ruch zwiększał lub zmniejszał prawdopodobieństwo wykonania tego ruchu w kolejnych krokach. To podejście miało jednak swoje wady, w szczególności opierało się wyłącznie na ostatniej nagrodzie, pomijając długoterminowe skutki.

W pierwszej próbie stworzyłem skrypt uczący się przez wzmocnianie pozytywnych akcji. Niestety, metoda ta działała tylko w oparciu o ostatnią nagrodę i nie uwzględniała długoterminowych efektów.

Poprawiłem ten problem, stosując strategię Epsilon-Greedy. Ta strategia uczy zarówno na podstawie natychmiastowych nagród, jak i poprzez eksploatację nowych możliwości. Dodatkowo, wprowadzono mechanizm cofania się do udanych wyborów z poprzednich iteracji, jeśli aktualna próba okazuje się mniej efektywna. W miarę upływu czasu agent mniej eksploruje, a bardziej polega na zdobytej wiedzy, zwiększając prawdopodobieństwo sukcesu.

3 Reinforcement Learning

W drugim podejściu zdecydowałem się skorzystać z biblioteki `stable_baselines3`, która umożliwia trenowanie agenta za pomocą zaawansowanych algorytmów, takich jak PPO (Proximal Policy Optimization) oraz A2C (Advantage Actor-Critic). Na początku wytrenowałem model przy użyciu algorytmu A2C. Chociaż rakietą była w stanie ładować, proces ten był stosunkowo powolny i nie zawsze kończył się sukcesem. Następnie wytrenowałem model za pomocą algorytmu PPO. Ten model uczył się znacznie szybciej i z większą precyzją, co przełożyło się na bardziej stabilne i efektywne ładowania.

3.1 Algorytm PPO

Algorytm Proximal Policy Optimization (PPO) jest jednym z najnowszych i najskuteczniejszych algorytmów uczenia przez wzmocnianie. Został opracowany przez OpenAI i jest szeroko stosowany ze względu na swoją prostotę, efektywność i stabilność.

Cel PPO

PPO dąży do poprawy polityki agenta (czyli strategii działania) w sposób stabilny i skuteczny, unikając dużych kroków aktualizacji, które mogą prowadzić do niestabilności.

Podejście do aktualizacji

Zamiast bezpośrednio maksymalizować funkcję celu, PPO wprowadza ograniczenia na wielkość zmian polityki, co pomaga utrzymać stabilność procesu uczenia.

Proces uczenia

1. **Generowanie doświadczeń:** Agent zbiera doświadczenia z interakcji ze środowiskiem, rejestrując stany, akcje, nagrody i stany następne.
2. **Obliczanie przewagi:** Na podstawie zebranych danych oblicza się wskaźnik przewagi dla każdej akcji.
3. **Aktualizacja polityki:** Aktualizacja parametrów polityki odbywa się przy użyciu funkcji celu PPO, co zapewnia, że zmiany są kontrolowane i stabilne.

Zalety PPO

- **Stabilność:** Dzięki mechanizmowi ograniczania zmian polityki, PPO jest bardziej stabilny w porównaniu do innych metod.
- **Prostota:** PPO jest prostszy do implementacji niż metody takie jak Trust Region Policy Optimization (TRPO), które również dążą do stabilnych aktualizacji.
- **Efektywność:** PPO efektywnie wykorzystuje dane, co czyni go wydajnym w praktyce.

PPO jest szeroko stosowany w różnych aplikacjach uczenia przez wzmocnienie ze względu na swoją efektywność i stabilność, co czyni go jednym z najpopularniejszych algorytmów w tej dziedzinie.

3.2 A2C

Algorytm Advantage Actor-Critic (A2C) jest jednym z popularnych algorytmów uczenia przez wzmocnienie. A2C łączy w sobie elementy metod policy gradient (aktora) i value-based (krytyka), co pozwala na efektywne trenowanie agenta.

Cel A2C

Celem A2C jest znalezienie optymalnej polityki działania (strategii), która maksymalizuje skumulowaną nagrodę otrzymywaną przez agenta w długim okresie. Algorytm osiąga to poprzez jednoczesne aktualizowanie dwóch sieci neuronowych: aktora i krytyka.

Aktor i Krytyk

- **Aktor:** Sieć neuronowa, która generuje politykę (policy), czyli prawdopodobieństwa wyboru konkretnych akcji w danym stanie. Aktor aktualizuje swoje parametry na podstawie wskaźnika przewagi (advantage).
- **Krytyk:** Sieć neuronowa, która estymuje funkcję wartości (value function), czyli oczekiwaną skumulowaną nagrodę dla danego stanu. Krytyk pomaga w ocenie jakości polityki generowanej przez aktora.

Wskaźnik przewagi

Wskaźnik przewagi (*advantage*) jest miarą jakości podjętej akcji w danym stanie w porównaniu do oczekiwań. Jest on obliczany jako różnica między rzeczywistą skumulowaną nagrodą a estymowaną wartością stanu:

$$A(s, a) = Q(s, a) - V(s)$$

gdzie:

- $Q(s, a)$ to wartość akcji a w stanie s (estymowana skumulowana nagroda).
- $V(s)$ to wartość stanu s (oczekiwana skumulowana nagroda).

Proces uczenia

1. **Generowanie doświadczeń:** Agent zbiera doświadczenia z interakcji ze środowiskiem, rejestrując stany, akcje, nagrody i stany następne.
2. **Obliczanie przewagi:** Na podstawie zebranych danych oblicza się wskaźnik przewagi dla każdej akcji.
3. **Aktualizacja aktora i krytyka:** Aktualizacja parametrów aktora i krytyka odbywa się przy użyciu odpowiednich funkcji celu, co zapewnia, że polityka i estymacje wartości stają się coraz dokładniejsze.

Zalety A2C

- **Efektywność:** A2C efektywnie łączy zalety metod policy gradient i value-based, co przyspiesza proces uczenia.
- **Stabilność:** Dzięki jednoczesnemu uczeniu aktora i krytyka, algorytm jest bardziej stabilny w porównaniu do metod opartych wyłącznie na policy gradient.
- **Elastyczność:** A2C może być stosowany w różnych środowiskach i problemach uczenia przez wzmacnianie.

A2C jest szeroko stosowany w aplikacjach uczenia przez wzmacnianie, oferując skuteczność i stabilność w trenowaniu agentów.

4 Kolejne kroki

Po udanym wytrenowaniu modelu na algorytmie PPO postanowiłem zmienić parametry środowiska, aby sprawdzić czy model jest w stanie nauczyć się lądować w innych warunkach. Zmieniłem m.in. grawitację, oraz dodałem wiatr. Obecnie modelowi czasami udaje się lądować, jednak nie jest to zbyt skuteczne. Postanowiłem go dotrenować na nowych warunkach. Trenowałem ten sam model zmieniając parametry środowiska. Po całkiem długim uczeniu model zaczął lądować w nowych warunkach. Obecnie model jest w stanie lądować w różnych warunkach (nawet ekstremalnych).

5 Algorytm genetyczny

Algorytmy genetyczne (AG) są technikami optymalizacyjnymi, które imitują procesy biologiczne. Populacja rozwiązań jest iteracyjnie ulepszana poprzez selekcję najlepszych osobników, krzyżowanie ich cech oraz wprowadzanie mutacji. W kontekście Lunar Lander, każde rozwiązanie (osobnik) reprezentuje sekwencję działań, które agent ma wykonać.

Kroki algorytmu genetycznego

1. **Inicjalizacja:** Generowanie początkowej populacji losowych rozwiązań (chromosomów).
2. **Selekcja:** Wybór najlepszych rozwiązań na podstawie funkcji oceny (fitness function).
3. **Krzyżowanie (Crossover):** Łączenie par wybranych rozwiązań w celu stworzenia nowych rozwiązań (potomstwa).
4. **Mutacja:** Wprowadzenie losowych zmian do nowych rozwiązań w celu zachowania różnorodności genetycznej.
5. **Ocena:** Obliczenie wartości funkcji oceny dla nowych rozwiązań.
6. **Zastąpienie:** Wybór najlepszych rozwiązań do nowej populacji.
7. **Iteracja:** Powtarzanie kroków 2-6 przez określoną liczbę generacji.

Funkcja Oceny (Fitness Function)

Funkcja oceny jest kluczowym elementem algorytmu genetycznego, ponieważ określa jakość poszczególnych rozwiązań. W naszym przypadku funkcja oceny ocenia sekwencje działań agenta na podstawie sumy nagród uzyskanych w trakcie jednego epizodu lądowania.

Kryteria oceny

- **Nagroda za lądowanie:** Agent otrzymuje pozytywne punkty za prawidłowe lądowanie.
- **Kara za dużą prędkość:** Jeżeli prędkość pionowa przekracza pewien próg, agent otrzymuje kare, aby zniechęcić do zbyt szybkiego opadania.
- **Kara za nachylenie:** Jeżeli nachylenie rakiety przekracza pewien kąt, agent otrzymuje kare, aby promować stabilne lądowanie.

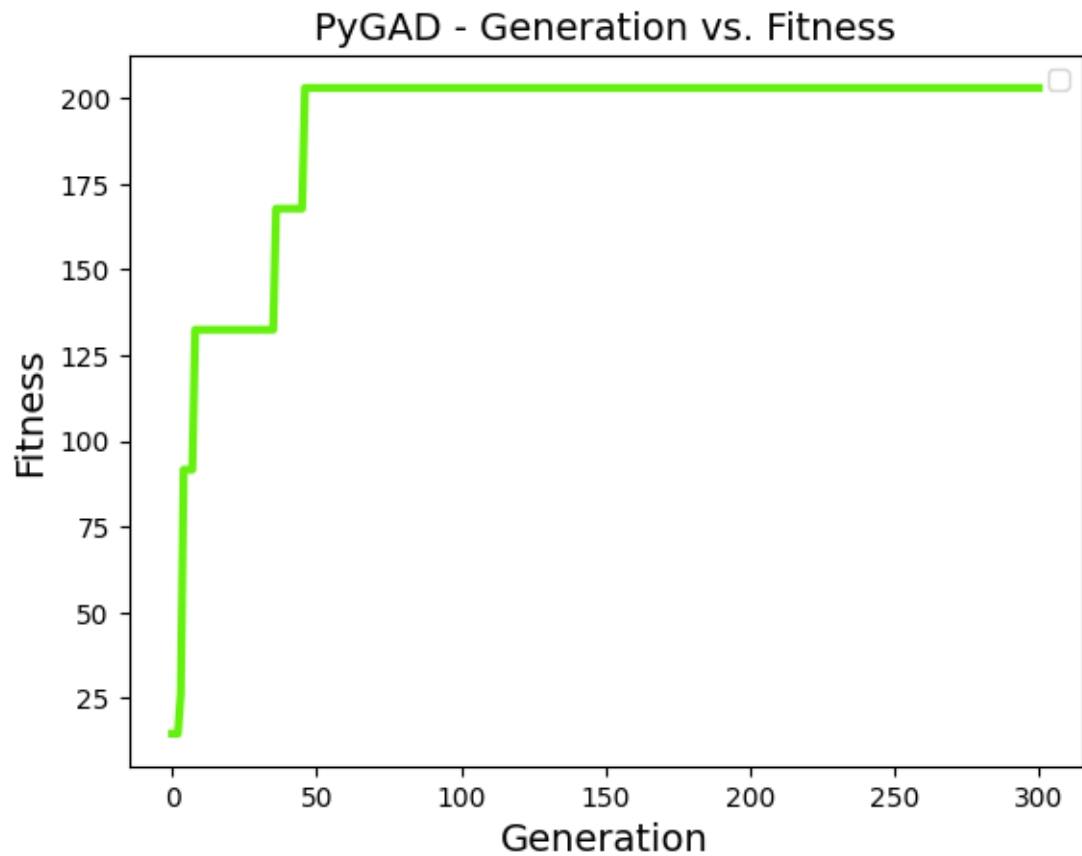
Łącząc powyższe kryteria, funkcja oceny weryfikuje, jak dobrze agent radzi sobie z zadaniem lądowania rakiety, nagradzając stabilne i kontrolowane lądowanie, jednocześnie karząc za ryzykowne manewry.

Parametry Algorytmu

Używamy następujących parametrów w naszym algorytmie genetycznym:

- **Populacja:** 100 chromosomów (rozwiązań).
- **Geny:** Każdy chromosom zawiera 200 genów, reprezentujących maksymalną liczbę ruchów.
- **Liczba pokoleń:** Algorytm jest uruchamiany przez 300 generacji.
- **Selekcja rodziców:** 50
- **Typ krzyżowania:** Używamy jednokrotnego krzyżowania (single-point crossover).
- **Mutacja:** Używamy typu mutacji *swap*, z domyślnym procentem genów podlegających mutacji.

Wyniki i Wizualizacja



6 Podsumowanie

Z przeprowadzonych eksperymentów wynika, że najlepsze wyniki osiągnięto stosując algorytmy uczenia przez wzmacnianie. Algorytmy te potrafiły efektywnie dostosować się do zmiennych warunków środowiskowych, takich jak zmiany grawitacji i wiatru. W przeciwieństwie do tego, algorytm genetyczny nie wykazał się równie dobrą adaptacyjnością. Definiował on stały zestaw ruchów, który nie był wystarczająco elastyczny, aby poradzić sobie z różnorodnością warunków początkowych, takich jak różne prędkości startowe czy zmienne nachylenie.

7 Materiały

- <https://www.youtube.com/watch?v=nRHjymV2PX8> YouTube - Reinforcement Learning
- <https://www.geeksforgeeks.org/epsilon-greedy-algorithm-in-reinforcement-learning> GeeksforGeeks - Epsilon-Greedy Algorithm in Reinforcement Learning
- https://gymnasium.farama.org/environments/box2d/lunar_lander Gymnasium - Lunar Lander Environment
- <https://github.com/sudharsan13296/Hands-On-Reinforcement-Learning-With-Python/blob/master/11>.
- <https://stable-baselines.readthedocs.io/en/master> Stable Baselines Documentation
- https://stable-baselines3.readthedocs.io/en/master/guide/rl_tips.html Stable Baselines3 - RL Tips
- <https://codecrucks.com/mamdani-fuzzy-inference-method-example> Codecrucks - Mamdani Fuzzy Inference Method Example