

# Zaawansowane Metody Inteligencji Obliczeniowej

## Lab 11: Ciągła Przestrzeń Akcji

Michał Kempka

Marek Wydmuch

20 maja 2021



Fundusze Europejskie  
Polska Cyfrowa



Rzeczpospolita  
Polska

Unia Europejska  
Europejski Fundusz  
Rozwoju Regionalnego



"Akademia Innowacyjnych Zastosowań Technologii Cyfrowych (AI Tech)",  
projekt finansowany ze środków Programu Operacyjnego Polska Cyfrowa POPC.03.02.00-00-0001/20

### 1 Actor-Critic dla ciągłej przestrzeni akcji

---

**Algorytm 1:** Pseudokod dla algorytmu One-step Actor-Critic

---

```
1 Inicjalizacja:  
2 Różniczkowalna polityka  $\pi(a|s, \theta)$   
3 Różniczkowalna funkcja wartości stanu  $v(s, w)$   
4 Zainicjalizowane wagi  $w \in \mathbb{R}_w^d$  i  $\theta \in \mathbb{R}_\theta^d$   
5 Zadana szybkość uczenia  $\alpha_\theta, \alpha_w \in \mathcal{R}^+$   
6 repeat  
7   if  $S$  nieustawiony lub terminalny then  
8     Rozpocznij nowy epizod  
9      $S \leftarrow S_0$   
10     $A \sim \pi(\cdot|S, \theta)$   
11    Wykonaj akcję  $A$ , zaobserwuj nagrodę  $R$  i następnik  $S'$   
12     $\delta \leftarrow R + \gamma v(S', w) - v(S, w)$   
13     $w \leftarrow w + \alpha_w \delta \nabla v(S, w)$   
14     $\theta \leftarrow \theta + \alpha_\theta \delta \nabla \ln \pi(A|S, \theta)$   
15     $S \leftarrow S'$   
16 until warunek stopu;
```

---

Na poprzednich zajęciach poznaliśmy metody gradientu polityki lub inaczej aproksymacji polityki, jedną z metod z tej rodziny jest algorytm Actor-Critic (Algorytm 1). Metody te oferują prosty sposób na radzenie sobie z dużymi przestrzeniami akcji, włączając to przestrzenie ciągłe, gdzie liczba możliwych akcji jest nieskończona. W wypadku ciągłej przestrzeni stanów zamiast uczyć się prawdopodobieństw dla każdej z akcji, możemy uczyć się parametrów rozkładu prawdopodobieństwa. Na przykład, przestrzenią akcji mogą być liczby rzeczywiste, a akcję będziemy wybierać z rozkładu normalnego. Funkcja gęstości rozkładu normalnego jest zdefiniowana w następujący sposób:

$$p(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)},$$

gdzie  $\mu$  i  $\sigma$  są średnią i odchyleniem standardowym rozkładu normalnego. Sparametryzowana polityka oparta na funkcji gęstości rozkładu normalnego jest więc następująca:

$$\pi(a|s, \theta) = \frac{1}{\sigma(s, \theta)\sqrt{2\pi}} e^{\left(-\frac{(a-\mu(s, \theta))^2}{2\sigma(s, \theta)^2}\right)},$$

gdzie  $\mu : \mathcal{S} \times \mathbb{R}^d \rightarrow \mathbb{R}, \sigma : \mathcal{S} \times \mathbb{R}^d \rightarrow \mathbb{R}^+$ . Dlatego podzielimy nasz wektor parametrów na dwie części, każda odpowiadająca za parametryzację jedynie jednej z tych funkcji:  $\theta = [\theta^\mu, \theta^\sigma]$ .

**Pytanie:** Reguła aktualizacji dla w algorytmie Actor-Critic jest następująca:

$$\theta_{t+1} = \theta_t + \alpha \delta_t \frac{\nabla \pi(A_t | S_t, \theta)}{\pi(A_t | S_t, \theta)} \quad (1)$$

$$= \theta_t + \alpha \delta_t \nabla \ln \pi(A_t | S_t, \theta). \quad (2)$$

Wyznacz regułę aktualizacji dla polityki opartej o sparametryzowaną funkcję gęstości rozkładu normalnego.

## 2 Deep Deterministic Policy Gradient (DDPG)

---

### Algorytm 2: Pseudokod dla algorytmu Deep Deterministic Policy Gradient

---

**1 Inicjalizacja:**

- 2 Różniczkowalne parametryczne funkcje rzeczywiste (sieci)  $Q(s, a | \theta^Q)$  i  $\mu(s | \theta^\mu)$
- 3 Powtórzone parametry dla 'target network':  $\theta^{Q'} \leftarrow \theta^Q, \theta^{\mu'} \leftarrow \theta^\mu$  Zainicjalizuj pamięć  $\mathcal{M} = \emptyset$
- 4 Parametr  $I$  określający co ile kroków uczymy się z pamięci
- 5 Parametr  $J$  określający ile doświadczeń wybieramy z pamięci
- 6 Parametr  $\tau \in (0, 1)$  odpowiadający za szybkość aktualizacji target network

7 steps  $\leftarrow 0$

**8 repeat**

9 steps  $\leftarrow$  steps + 1

10 **if**  $S$  nieustawiony lub terminalny **then**

11     Rozpocznij nowy epizod

12      $S \leftarrow S_0$

13     Wybierz  $A$  według dowolnej polityki (np. zgodnej z  $\mu$  poszerzonej o losowość)

14     Wykonaj akcję  $A$ , zaobserwuj nagrodę  $R$  i następnik  $S'$

15      $\mathcal{M} \leftarrow \mathcal{M} \cup \langle S, A, R, S' \rangle$

16      $S \leftarrow S'$

17     **if** steps %  $I = 0$  **then**

18          $g \leftarrow 0$

19         **for**  $j \leftarrow 1, \dots, J$  **do**

20              $\langle S, A, R, S' \rangle \leftarrow$  losowo wybrane z  $\mathcal{M}$

21              $g^Q \leftarrow g^Q + \frac{1}{2} \nabla_{\theta^Q} \left( Q(S, A | \theta^Q) - (R + \gamma Q'(S', \mu'(S' | \theta^{\mu'})) | \theta^{Q'}) \right)^2$

22              $g^\mu \leftarrow g^\mu + \nabla_{\theta^\mu} Q(S, \mu(s | \theta^\mu) | \theta^Q)$

23              $\theta^Q \leftarrow \theta^Q - \frac{1}{J} \alpha g^Q$

24              $\theta^\mu \leftarrow \theta^\mu + \frac{1}{J} \alpha g^\mu$

25              $\theta^{Q'} \leftarrow \tau \theta^Q + (1 - \tau) \theta^{Q'}$

26              $\theta^{\mu'} \leftarrow \tau \theta^\mu + (1 - \tau) \theta^{\mu'}$

27 **until** warunek stopu;

---

Deep Deterministic Policy Gradient (DDPG) [1] to podejście, które łączy metodę Actor-Critic z metodą Deep Q-Learning. Pozwala nam użyć ciągłych akcji, jest **off-policy**, **model-free** i w naturalny sposób może korzystać ze wszelkich (większości?) rozszerzeń metody DQN. W praktyce możemy nawet zapomnieć o fakcie powiązania DDPG z metodami z rodziny policy gradient i organicznie wymyślić ją wychodząc od DQN.

### 2.1 Ciągłe akcje w DQN

W normalnym DQN, wejście dla sieci jest stanem, a wyjście wektorem  $Q$  długości przestrzeni akcji. W przypadku DDPG jest to niemożliwe ze względu na fakt, że zbiór akcji jest niepoliczalny (lub prawdopodobnie zaporowo duży w zdyskretyzowaniu). Możemy więc użyć zatem naszej ciągłej akcji (wektora liczb rzeczywistych) jako części **wejścia** do sieci. Zamiast  $Q_{DQN} : \mathcal{S} \rightarrow \mathbb{R}^{|\mathcal{A}|}$  uzyskujemy  $Q_{DDPG} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ .

Niestety przy takiej formulacji pojawia się pewien problem: jak wybrać akcję maksymalizującą  $Q$ ? Stwórzmy funkcję  $\mu(S) : \mathcal{S} \rightarrow \mathbb{R}^d$  (gdzie  $d$  to wymiarowość akcji), która zwraca nam **najlepszą akcję** dla danego stanu

(preferencyjnie wyrocznia). Możemy teraz użyć  $\mu(S')$  w klasycznej aktualizacji Q-learningu pozbywając się maxa:  $\max_{a' \in A} Q(S', a') \rightarrow Q(S', \mu(S'|\theta^\mu))$ .

Założenie, że  $\mu$  jest wyrocznią nie jest niestety bardzo pomocne w praktyce więc określmy, że  $\mu$  będzie funkcją parametryczną z parametrami  $\theta^\mu$ . Z uwagi na fakt, że **zwracanie najlepszej akcji** jest tożsame z **maksymalizacją Q** możemy sformułować podproblem optymalizacyjny maksymalizacji Q ze względu na  $\theta^\mu$  (średnio dla wszystkich możliwych stanów):

$$\max_{\theta^\mu} \mathbb{E}_{s \in S} [Q(s, \mu(s|\theta^\mu))] \quad (3)$$

Powyższy podproblem możemy rozwiązać (rozwiązywać) na różne sposoby, Ale naturalnym, prostym i przyjemnym rozwiązaniem jest użycie SGD z przeciwnym znakiem - **gradient ascent** i przeprowadzanie optymalizacji równoległe z klasyczną aktualizacją funkcji Q. Zakładając, że funkcja Q jest także parametryzowana, powiedzmy przez  $\theta^Q$ , przy napotkaniu pojedynczej krotki  $\{S, A, S', R\}$  możemy zastosować dwie do pewnego stopnia niezależne aktualizacje:

$$\theta^Q \leftarrow \theta^Q - \alpha \frac{1}{2} \nabla_{\theta^Q} \left( Q(S, A|\theta^Q) - (R + \gamma Q(S', \mu(S'|\theta^\mu)|\theta^Q)) \right)^2 \quad (4)$$

$$\theta^\mu \leftarrow \theta^\mu + \alpha \nabla_{\theta^\mu} Q(S, \mu(S|\theta^\mu)|\theta^Q) \quad (5)$$

Pseudokod 2 prezentuje ideę DDPG bardziej szczegółowo i z użyciem standardowych ulepszeń z DQN (experience replay, target network)

## 2.2 Architektura sieci i parametry

Jak widać z równań 4 mamy dwa zestawy parametrów, które mogą być użyte do zaimplementowania dwóch osobnych sieci, które będą aktualizowane niezależnie. Bardzo popularne jednak jest, by część sieci **odpowiedzialna za przetwarzanie stanów** była współdzielona.

## 2.3 Nawiązanie do Actor-Critic

Jak wspomniano wcześniej, DDPG możemy traktować jako metodę z rodziny Actor-Critic gdzie krytyk uczy się  $Q(S, a)$ , a aktor polityki  $\mu(s)$ .

## 2.4 Eksploracja

Z uwagi na fakt, że DDPG, podobnie jak Q-learning jest metodą off-policy, nasza polityka może być teoretycznie dowolna. Typowym jednak jest by stosować eksplorację poprzez wprowadzenie szumu Gaussowskiego (rozkład normalny) skupionego na 0 z progresywnie malejącą wariancją. Można jednak zastosować inny model eksploracji bez większych teoretycznych reperkusji.

## 2.5 Powolna aktualizacja sieci

Jak widać w pseudokodzie 2 (ostatnie linijki) w DDPG mamy podobnie jak w DQN do czynienia z 'zamrożonymi' sieciami (target network), które w tym przypadku są aktualizowane na bieżąco, lecz tylko częściowo. Nie jest to kluczowa część algorytmu, lecz warto zwrócić uwagę na istnienie takiej alternatywy.

## 2.6 Rozkłady akcji

Do tej pory, mówiąc o akcjach ciągłych, domniemywaliśmy, że chodzi o wektory liczb rzeczywistych. W ogólności nie musi tak być (choć jest to raczej rzadziej poruszany temat). Na akcje mogą być narzucone najróżniejsze ograniczenia np. mają być dodatnie, z pewnego zakresu  $[a, b]$ , tylko liczby naturalne. W tym przypadku należy użyć odpowiednich funkcji aktywacji (np. sigmoid, softplus) dla wyjść sieci i rozkładów prawdopodobieństw z odpowiednimi funkcjami gęstości (np. ucięty rozkład normalny, rozkład Bernoulliego)

## Literatura

- [1] Lillicrap, T. P., Hunt, J. J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., Silver, D., and Wierstra, D. (2016). Continuous control with deep reinforcement learning. In Bengio, Y. and LeCun, Y., editors, *ICLR*.
- [2] Russell, S. and Norvig, P. (2010). *Artificial Intelligence: A Modern Approach*. Prentice Hall, third edition.
- [3] Sutton, R. S. and Barto, A. G. (2018). *Reinforcement Learning: An Introduction*. The MIT Press, second edition.



**Fundusze  
Europejskie**  
Polska Cyfrowa



**Rzeczpospolita  
Polska**

**Unia Europejska**  
Europejski Fundusz  
Rozwoju Regionalnego



"Akademia Innowacyjnych Zastosowań Technologii Cyfrowych (AI Tech)",  
projekt finansowany ze środków Programu Operacyjnego Polska Cyfrowa POPC.03.02.00-00-0001/20