

Zaawansowane Metody Inteligencji Obliczeniowej

Wykład 3: Procesy Decyzyjne Markowa

Michał Kempka Marek Wydmuch Bartosz Wieloch

14 marca 2022



**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

Plan wykładu

1 Procesy Decyzyjne Markowa

2 Programowanie dynamiczne

Agent i środowisko

- **Agent** — podejmuje decyzje (wykonuje **akcje**) i się uczy maksymalizować nagrody
- **Środowisko** — reaguje na akcje zmieniając swój **stan** oraz zwracając **nagrody**

$$S_0, A_0, R_1, S_1, A_1, R_2, S_2, A_2, R_3, S_3, A_3, \dots$$

- po wykonaniu akcji A_t agent dostaje nagrodę R_{t+1} wraz z nowym stanem S_{t+1} (różne konwencje w literaturze!)

Procesy Decyzyjne Markowa (MDP)

Procesy Decyzyjne Markowa — będziemy stosowali skrót **MDP** (*ang. **M**arkov **D**ecision **P**rocesses*)

Co to jest Proces Decyzyjny Markowa?

\mathcal{S} — zbiór stanów

\mathcal{A} — zbiór akcji (gdy zależny od stanu: $\mathcal{A}(s)$)

\mathcal{R} — zbiór nagród

$p(s', r|s, a)$ — funkcja przejść (definiuje dynamikę MDP)

- $p(s', r|s, a) = P(S_t = s', R_t = r | S_{t-1} = s, A_{t-1} = a)$
Jest to zwykła, **deterministyczna** funkcja 4 argumentowa.
- W MDP funkcja p całkowicie określa dynamikę środowiska:
prawdopodobieństwa S_t oraz R_t zależą tylko i wyłącznie od poprzedzającego stanu S_{t-1} i akcji A_{t-1} (własność Markowa)

MDP

W MDP na podstawie 4-argumentowej funkcji p możemy wyliczyć wszystko co dotyczy środowiska, np.:

- funkcja przejść stanów

$$p(s'|s, a) = P(S_t = s' | S_{t-1} = s, A_{t-1} = a) = \sum_{r \in \mathcal{R}} p(s', r | s, a)$$

- spodziewana nagroda dla pary stan-akcja

$$r(s, a) = \mathbb{E} [R_t | S_{t-1} = s, A_{t-1} = a] = \sum_{r \in \mathcal{R}} r \sum_{s' \in \mathcal{S}} p(s', r | s, a)$$

- spodziewana nagroda dla pary stan-akcja-stan

$$r(s, a, s') = \mathbb{E} [R_t | S_{t-1} = s, A_{t-1} = a, S_t = s'] = \sum_{r \in \mathcal{R}} r \frac{p(s', r | s, a)}{p(s' | s, a)}$$

Cele i nagrody

Celem agenta jest **maksymalizacja oczekiwanej całkowitej** nagrody którą dostanie (maksymalizacja nie tylko natychmiastowej nagrody ale **skumulowanej w długim okresie**).

Przykłady:

- robot uczący się chodzić: w każdym kroku nagroda proporcjonalna do przebytego dystansu
- robot wychodzący z labiryntu: -1 za każdy krok zanim znajdzie wyjście
- robot zbierający śmieci: $+1$ za każdy znaleziony i podniesiony śmieć (-1 za uderzenie w przeszkody itd.)

Nagroda **nie powinna** sugerować jak osiągnąć pożądany cel (np. w grze w szachy agent nie powinien być nagradzany za zabicie przeciwnika itp.)

Spodziewany zysk

Formalnie chcemy maksymalizować oczekiwany zysk G_t od chwili t .

G_t to funkcja sekwencji nagród $R_{t+1}, R_{t+2}, R_{t+3}, \dots$

- Suma nagród:

$$G_t = R_{t+1} + R_{t+2} + R_{t+3} + \dots + R_T,$$

gdzie T to ostatni krok.

Ma sens gdy w problemie naturalnie występuje „ostatni” krok (np. koniec gry).

Interakcje agent–środowisko można podzielić w takim wypadku na niezależne epizody.

- Suma **zdyskontowanych** nagród:

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1},$$

gdzie $0 \leq \gamma \leq 1$ jest współczynnikiem dyskontowym.

- ▶ $\gamma = 0$ — agent jest „krótkowzroczny”
- ▶ $\gamma \rightarrow 1$ — agent staje się „dalekowzroczny”: przyszłe nagrody liczą się coraz bardziej

Wzór rekurencyjny

Dla zdyskontowanych nagród mamy:

$$\begin{aligned} G_t &= R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4} + \dots \\ &= R_{t+1} + \gamma (R_{t+2} + \gamma R_{t+3} + \gamma^2 R_{t+4} + \dots) \\ &= R_{t+1} + \gamma G_{t+1} \end{aligned}$$

Polityka

Polityka π — funkcja mapująca stan na prawdopodobieństwo wybrania każdej możliwej akcji

$$\pi(a|s) = P(A_t = a | S_t = s)$$

π jest zwykłą funkcją dwuargumentową (a „|” tylko przypomina, że definiuje ona rozkład prawdopodobieństwa).

Metody uczenia ze wzmocnieniem określają jak polityka agenta zmienia się wraz z jego doświadczeniem.

Funkcja wartości stanu

$v_\pi(s)$ — funkcja **wartości stanu** (inaczej: użyteczności stanu) dla polityki π określa **oczekiwany zysk** gdy:

- 1 rozpoczynamy w stanie s
- 2 dalej podążamy za polityką π

$$\begin{aligned} v_\pi(s) &= \mathbb{E}_\pi [G_t | S_t = s] \\ &= \mathbb{E}_\pi \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s \right] \end{aligned}$$

$\mathbb{E}_\pi [\cdot]$ — oczekiwana wartość zmiennej losowej gdy agent stosuje politykę π

Funkcja wartości akcji

$q_\pi(s, a)$ — funkcja wartości akcji a podjętej w stanie s dla polityki π określa oczekiwany zysk gdy:

- 1 rozpoczynamy w stanie s
- 2 wykonujemy akcję a
- 3 dalej stosujemy politykę π

$$\begin{aligned} q_\pi(s, a) &= \mathbb{E}_\pi [G_t | S_t = s, A_t = a] \\ &= \mathbb{E}_\pi \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s, A_t = a \right] \end{aligned}$$

Równanie Bellmana

Zależność pomiędzy wartością stanu s a wartością możliwych następnych stanów:

$$\begin{aligned}v_{\pi}(s) &= \mathbb{E}_{\pi} [G_t | S_t = s] \\&= \mathbb{E}_{\pi} [R_{t+1} + \gamma G_{t+1} | S_t = s] \\&= \sum_a \pi(a|s) \sum_{s'} \sum_r p(s', r | s, a) [r + \gamma \mathbb{E}_{\pi} [G_{t+1} | S_{t+1} = s']] \\&= \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) [r + \gamma v_{\pi}(s')]\end{aligned}$$

Inaczej:

- suma $[r + \gamma v_{\pi}(s')]$
- po wszystkich możliwych a , s' i r
- ważona prawdopodobieństwem $\pi(a|s)p(s', r | s, a)$

Równanie Bellmana

$$v_{\pi}(s) = \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a) [r + \gamma v_{\pi}(s')]$$

$$v_{\pi}(s) = \sum_{a,s',r} \pi(a|s) p(s',r|s,a) [r + \gamma v_{\pi}(s')]$$

Z równania Bellmana można wyliczyć v_{π} — mamy układ równań (dla każdego $s \in \mathcal{S}$)

Równanie Bellmana

$$v_{\pi}(s) = \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a) [r + \gamma v_{\pi}(s')]$$

$$v_{\pi}(s) = \sum_{a,s',r} \pi(a|s) p(s',r|s,a) [r + \gamma v_{\pi}(s')]$$

Z równania Bellmana można wyliczyć v_{π} — mamy układ równań (dla każdego $s \in \mathcal{S}$)

Pytanie

Czy łatwo jest wyliczyć $v_{\pi}(s)$ dla każdego $s \in \mathcal{S}$ korzystając z równania Bellmana?

Optymalna polityka — funkcja wartości i akcji

Dla MDP można zdefiniować optymalną politykę — jest to polityka π maksymalizująca $v_\pi(s)$ dla wszystkich stanów s .

Optymalna funkcja wartości:

$$v_*(s) = \max_{\pi} v_\pi(s)$$

dla każdego $s \in \mathcal{S}$

Optymalna funkcja akcji:

$$q_*(s, a) = \max_{\pi} q_\pi(s, a)$$

dla każdego $s \in \mathcal{S}$ oraz $a \in \mathcal{A}$

Zależność optymalnej funkcji akcji od funkcji wartości:

$$q_*(s, a) = \mathbb{E} [R_{t+1} + \gamma v_*(S_{t+1}) | S_t = s, A_t = a]$$

Optymalna polityka — równanie Bellmana

Równanie Bellmana dla optymalnej funkcji wartości stanu — bez odwołania do żadnej specyficznej polityki π .

$$\begin{aligned}v_*(s) &= \max_{a \in \mathcal{A}(s)} q_{\pi_*}(s, a) \\&= \max_a \mathbb{E}_{\pi_*} [G_t | S_t = s, A_t = a] \\&= \max_a \mathbb{E}_{\pi_*} [R_{t+1} + \gamma G_{t+1} | S_t = s, A_t = a] \\&= \max_a \mathbb{E} [R_{t+1} + \gamma v_*(S_{t+1}) | S_t = s, A_t = a] \\&= \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma v_*(s')]\end{aligned}$$

Optymalna polityka — równanie Bellmana

Równanie Bellmana dla optymalnej funkcji wartości stanu — bez odwołania do żadnej specyficznej polityki π .

$$\begin{aligned}v_*(s) &= \max_{a \in \mathcal{A}(s)} q_{\pi_*}(s, a) \\&= \max_a \mathbb{E}_{\pi_*} [G_t | S_t = s, A_t = a] \\&= \max_a \mathbb{E}_{\pi_*} [R_{t+1} + \gamma G_{t+1} | S_t = s, A_t = a] \\&= \max_a \mathbb{E} [R_{t+1} + \gamma v_*(S_{t+1}) | S_t = s, A_t = a] \\&= \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma v_*(s')]\end{aligned}$$

Pytanie

Czy łatwo jest wyliczyć $v_*(s)$ dla każdego $s \in \mathcal{S}$ korzystając z powyższego równania Bellmana?

Optymalna polityka – równanie Bellmana

Równanie Bellmana dla optymalnej funkcji wartości akcji:

$$\begin{aligned} q_*(s, a) &= \mathbb{E} \left[R_{t+1} + \gamma \max_{a'} q_*(S_{t+1}, a') \mid S_t = s, A_t = a \right] \\ &= \sum_{s', r} p(s', r \mid s, a) \left[r + \gamma \max_{a'} q_*(s', a') \right] \end{aligned}$$

Optymalna polityka — rozwiązanie

Mając daną funkcję wartości stanu $v_*(s)$ optymalna polityka to taka która przypisuje niezerowe prawdopodobieństwo tylko i wyłącznie dla akcji które mają maksymalną **oczekiwaną** wartość:

$$\pi_*(s|a) > 0 \Rightarrow a \in \arg \max_a \sum_{s', r} p(s', r|s, a) v_*(s')$$

Mając daną funkcję wartości akcji $q_*(s, a)$ optymalna polityka to:

$$\pi_*(s|a) > 0 \Rightarrow a \in \arg \max_a q_*(s, a)$$

Optymalna polityka — rozwiązanie

Mając daną funkcję wartości stanu $v_*(s)$ optymalna polityka to taka która przypisuje niezerowe prawdopodobieństwo tylko i wyłącznie dla akcji które mają maksymalną **oczekiwaną** wartość:

$$\pi_*(s|a) > 0 \Rightarrow a \in \arg \max_a \sum_{s',r} p(s',r|s,a) v_*(s')$$

Mając daną funkcję wartości akcji $q_*(s,a)$ optymalna polityka to:

$$\pi_*(s|a) > 0 \Rightarrow a \in \arg \max_a q_*(s,a)$$

Pytanie

Wyznaczenie optymalnej polityki **wymaga modelu świata** dla metody opartej na funkcji wartości stanu, wartości akcji czy dla obu?

Plan wykładu

- 1 Procesy Decyzyjne Markowa
- 2 Programowanie dynamiczne

Ewaluacja polityki

Ewaluacja polityki (*ang. Policy Evaluation*) — proces wyliczania funkcji wartości stanu $v_\pi(s)$ dla konkretnej polityki π .

$$\begin{aligned} v_\pi(s) &= \mathbb{E}_\pi [G_t | S_t = s] \\ &= \mathbb{E}_\pi [R_{t+1} + \gamma G_{t+1} | S_t = s] \\ &= \mathbb{E}_\pi [R_{t+1} + \gamma v_\pi(S_{t+1}) | S_t = s] \\ &= \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) [r + \gamma v_\pi(s')] \end{aligned}$$

- gdy dynamika problemu całkowicie znana (MDP) to mamy układ $|S|$ równań z $|S|$ niewiadomymi
- rozwiązanie jest proste (aczkolwiek może być pracochłonne!) — mamy gotowe metody rozwiązywania układów równań
- alternatywa: **metody iteracyjne**

Ewaluacja polityki — metoda iteracyjna

W metodzie iteracyjnej rozważamy **sekwencję przybliżeń** funkcji wartości: v_0, v_1, v_2, \dots

- początkowe przybliżenie wybrane arbitralnie (np. same zera)
- każde kolejne przybliżenie uzyskane z równania Bellmana:

$$\begin{aligned} v_{k+1}(s) &= \mathbb{E}_{\pi} [R_{t+1} + \gamma v_k(S_{t+1}) | S_t = s] \\ &= \sum_a \pi(a|s) \sum_{s', r} p(s', r|s, a) [r + \gamma v_k(s')] \end{aligned}$$

dla wszystkich $s \in \mathcal{S}$

- sekwencja v_k zbiega do v_{π} dla $k \rightarrow \infty$
- implementacja powyższej aktualizacji wymaga dwóch tablic: dla $v_{k+1}(s)$ oraz dla $v_k(s)$
- aktualizacja „w miejscu” — mamy tylko jedną tablicę, aktualizacja wartości dla stanu s może zależeć zarówno od starych jak i już zaktualizowanych wartości (również zbiega do v_{π})

Algorytm 1: Pseudokod dla algorytmu ewaluacji polityki (ang. Policy Evaluation)

1 Argumenty:

2 Polityka $\pi(s)$ dla wszystkich $s \in \mathcal{S}$

3 $\Delta_{min} \leftarrow$ mała dodatnia liczba

4 Inicjalizacja:

5 Dowolnie zainicjalizowane $V(s)$ (np. $V(s) = 0$) – estymat wartości stanów $v(s)$

6 repeat

7 $\Delta \leftarrow 0$

8 **for** $s \in \mathcal{S}$ **do**

9 $v \leftarrow V(s)$

10 $V(s) \leftarrow \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a)[r + \gamma V(s')]$

11 $\Delta \leftarrow \max(\Delta, |v - V(s)|)$

12 **until** $\Delta < \Delta_{min}$;

Ulepszenie polityki

Celem ewaluacji polityki jest znalezienie lepszej polityki :)

- wiemy jak dobrze jest postępować zgodnie z daną polityką π
- czy warto w stanie s użyć innej akcji a i dalej podążać już zgodnie z π ?
- wartość takiej akcji jest równa:

$$\begin{aligned} q_{\pi}(s, a) &= \mathbb{E} [R_{t+1} + \gamma v_{\pi}(S_{t+1}) | S_t = s, A_t = a] \\ &= \sum_{s', r} p(s', r | s, a) [r + \gamma v_{\pi}(s')] \end{aligned}$$

- jeśli $q_{\pi}(s, a) > v_{\pi}(s)$
to opłaca się zmienić politykę dla stanu s !
(mamy gwarancję poprawy polityki)

Algorytm iteracji polityki (ang. *Policy Iteration*)

Algorytm 2: Pseudokod dla algorytmu iteracji polityki (ang. Policy Iteration)

```
1  1. Inicjalizacja:
2  Dowolnie zainicjalizowane  $V(s)$  (np.  $V(s) = 0$ ) i  $\pi(s)$  dla wszystkich  $s \in \mathcal{S}$ 
3   $\Delta_{min} \leftarrow$  mała dodatnia liczba

4  2. Ewaluacja polityki:
5  repeat
6       $\Delta \leftarrow 0$ 
7      for  $s \in \mathcal{S}$  do
8           $v \leftarrow V(s)$ 
9           $V(s) \leftarrow \sum_{s',r} p(s', r|s, a)[r + \gamma V(s')]$ 
10          $\Delta \leftarrow \max(\Delta, |v - V(s)|)$ 
11 until  $\Delta < \Delta_{min};$ 

12 3. Polepszenie polityki:
13  $Stop \leftarrow True$ 
14 for  $s \in \mathcal{S}$  do
15      $a \leftarrow \pi(s)$ 
16      $\pi(s) \leftarrow$ 
         $\arg \max_a \sum_{s',r} p(s', r|s, \pi(s))[r + \gamma V(s')]$ 
17     if  $a \neq \pi(s)$  then  $Stop \leftarrow False$ 
18 if  $Stop$  then return  $\pi$  else goto 2
```

Algorytm iteracji polityki — uwagi

- Jeśli θ będzie zbyt duże to procedura ewaluacji polityki (krok 2) zakończy się zbyt szybko
 - ▶ nawet bardzo niewielka poprawa dokładności $V(s)$ może spowodować wybranie przez $\arg \max$ innej akcji w procedurze ulepszenia polityki (krok 3)
 - ▶ może to skutkować ustabilizowaniem polityki π na suboptymalnym rozwiązaniu
- W procedurze ulepszenia polityki, polityka staje się stabilna gdy **zbiór** najlepszych akcji (maksymalizujących wartość akcji) się nie zmienia.
Jeśli $\arg \max$ zwraca tylko jedną arbitralną akcję może dojść do niekończącego przełączania się między dwoma równie dobrymi akcjami.
- W procedurze ewaluacji polityki zaczynamy od aktualnego oszacowania V z poprzedniego kroku — szybsza zbieżność.

Algorytm iteracji wartości (*ang. Value Iteration*)

Algorytm iteracji wartości to specyficzny przypadek algorytmu iteracji polityki:

- **jedna** iteracja ewaluacji polityki
- ulepszenie polityki

W praktyce nie wyznacza się jawnie polityki w każdej iteracji — powyższe dwa kroki łączy się w jeden, a politykę wyznacza się dopiero na końcu.

Algorytm iteracji wartości

Algorytm 3: Pseudokod dla algorytmu **iteracji wartości** (ang. Value Iteration)

```
1 inicjalizacja:  
2 Dłownie zainicjalizowane  $V(s)$  dla wszystkich  $s \in \mathcal{S}$   
3  $\Delta_{min} \leftarrow$  mała dodatnia liczba  
4 repeat  
5      $\Delta \leftarrow 0$   
6     for  $s_i \in \mathcal{S}$  do  
7          $v \leftarrow V(s)$   
8          $V(s) \leftarrow \max_a \sum_{s',r} p((s',r|s,a)[r + \gamma V(s')]$   
9          $\Delta \leftarrow \max(\Delta, |v - V(s)|)$   
10 until  $\Delta < \Delta_{min}$ ;  
11 for  $s_i \in \mathcal{S}$  do  
12      $\pi(s) \leftarrow \arg \max_a \sum_{s',r} p((s',r|s,a)[r + \gamma V(s')]$   
13 return  $\pi$ 
```

Zmodyfikowany algorytm iteracji polityki

Wadą poprzednio omawianego algorytmu iteracji polityki jest koszt obliczeniowy:

- w głównej pętli występuje procedura ewaluacji polityki, która...
- też jest metodą iteracyjną.

Pomysł: zastosować przybliżoną wartość polityki $V(s)$ — wykonujemy k kroków algorytmu iteracji wartości

Zmodyfikowany algorytm iteracji polityki

Algorytm 4: Pseudokod dla algorytmu iteracji polityki z k

```
1 1. Inicjalizacja:  
2 Dowolnie zainicjalizowane  $V(s)$  (np.  $V(s) = 0$ ) i  $\pi(s)$  dla wszystkich  $s \in \mathcal{S}$   
3  $\Delta_{min} \leftarrow$  mała dodatnia liczba;  $k \leftarrow$  dodatnia liczba całkowita  
4 2. Ewaluacja polityki:  
5 for  $i = 0; i < k; i \leftarrow i + 1$  do  
6   for  $s \in \mathcal{S}$  do  
7      $V(s) \leftarrow \sum_{s', r} p(s', r | s, a)[r + \gamma V(s')]$   
8 3. Polepszenie polityki:  
9  $Stop \leftarrow True$   
10 for  $s \in \mathcal{S}$  do  
11    $a \leftarrow \pi(s)$   
12    $\pi(s) \leftarrow$   
        $\arg \max_a \sum_{s', r} p(s', r | s, \pi(s))[r + \gamma V(s')]$   
13   if  $a \neq \pi(s)$  then  $Stop \leftarrow False$   
14 if  $Stop$  then return  $\pi$  else goto 2
```

Pytanie

Co się stanie gdy $k = 1$?

Asynchroniczne programowanie dynamiczne

- Wadą omówionych metod jest to, że wymagają wykonywania pewnych operacji po wszystkich stanach.
- Dla dużej liczby stanów nawet pojedyncza iteracja bardzo kosztowna.
- Idea: nie przechodzić systematycznie po wszystkich stanach w każdej iteracji.
 - ▶ Aktualizacja wartości stanów w dowolnej kolejności
 - ▶ Dla niektórych stanów aktualizacje częściej, dla innych – rzadziej.
- Dla zbieżności nadal konieczna jest aktualizacja wartości dla wszystkich stanów!
- Ale możemy częściej aktualizować wartości dla stanów które interesują agenta.

Bibliografia

- [1] Russell, S. and Norvig, P. (2010). *Artificial Intelligence: A Modern Approach*. Prentice Hall, third edition.
 - [2] 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