

## STRATEGIA OPTYMALNA

**Karol Działowski**

nr albumu: 39259

przedmiot: Uczenie ze wzmocnieniem

Szczecin, 24 grudnia 2020

### Spis treści

<b>1 Cel laboratorium</b>	<b>1</b>
<b>2 Zadanie 3</b>	<b>1</b>
2.1 Równania Bellmana . . . . .	2
2.2 Wyznaczone wartości strategii . . . . .	3
2.3 Optymalna strategia gry . . . . .	4
2.4 Kod źródłowy . . . . .	4
<b>3 Zadanie 4</b>	<b>7</b>
3.1 Założenia . . . . .	7
3.2 Wyniki . . . . .	7
3.3 Kod źródłowy . . . . .	8

### 1 Cel laboratorium

Celem laboratorium nr 2 było wykorzystanie równań Bellmana do wyznaczania strategii optymalnej.

### 2 Zadanie 3

Przedmiotem rozważań jest gra w zapalki. Gracze mają do dyspozycji pewną liczbę (np.  $N = 10$ ) zapalek. W grze bierze udział dwóch graczy, którzy na przemian zabierają 1, 2 lub 3

zapałki. Przegrywa ten gracz, który zabiera zapałki jako ostatni.

Należało wyznaczyć funkcję wartości dla strategii:  $\pi_1(x) = -1$ ,  $\pi_2(x) = -2$ ,  $\pi_3(x) = -3$  oraz określić, która z danych strategii jest najlepsza.

Ponadto wyznaczono optymalną strategię gry, zakładając, że nasz przeciwnik gra w sposób losowy.

Oczekiwana wartość wzmocnienia:

$$R_{xy}^a = \begin{cases} 1 & \text{dla } (x = 2, a = -1, y = 0) \text{ lub } (x = 3, a = -2, y = 0) \text{ lub } (x = 4, a = -3, y = 0) \\ 0 & \text{w innych przypadkach} \end{cases} \quad (1)$$

Prawdopodobieństwa przejścia rozpisano dla każdego stanu:

Stan 0:

$$P_{00}^{-1} = 1 \quad P_{00}^{-2} = 1 \quad P_{00}^{-3} = 1$$

Stan 1:

$$P_{10}^{-1} = 1 \quad P_{10}^{-2} = 1 \quad P_{10}^{-3} = 1$$

Stan 2:

$$P_{20}^{-1} = 1 \quad P_{20}^{-2} = 1 \quad P_{20}^{-3} = 1$$

Stan 3:

$$P_{30}^{-1} = \frac{2}{3} \quad P_{30}^{-2} = 1 \quad P_{30}^{-3} = 1$$

$$P_{31}^{-1} = \frac{1}{3}$$

Stan 4:

$$P_{40}^{-1} = \frac{1}{3} \quad P_{40}^{-2} = \frac{2}{3} \quad P_{40}^{-3} = 1$$

$$P_{41}^{-1} = \frac{1}{3} \quad P_{41}^{-2} = \frac{1}{3}$$

$$P_{42}^{-1} = \frac{1}{3}$$

Stan 5:

$$P_{51}^{-1} = \frac{1}{3} \quad P_{50}^{-2} = \frac{1}{3} \quad P_{50}^{-3} = \frac{2}{3}$$

$$P_{52}^{-1} = \frac{1}{3} \quad P_{51}^{-2} = \frac{1}{3} \quad P_{51}^{-3} = \frac{1}{3}$$

$$P_{53}^{-1} = \frac{1}{3} \quad P_{52}^{-2} = \frac{1}{3}$$

Stan 6:

$$P_{62}^{-1} = \frac{1}{3} \quad P_{61}^{-2} = \frac{1}{3} \quad P_{60}^{-3} = \frac{1}{3}$$

$$P_{63}^{-1} = \frac{1}{3} \quad P_{62}^{-2} = \frac{1}{3} \quad P_{61}^{-3} = \frac{1}{3}$$

$$P_{64}^{-1} = \frac{1}{3} \quad P_{63}^{-2} = \frac{1}{3} \quad P_{62}^{-3} = \frac{1}{3}$$

Stan 7:

$$P_{73}^{-1} = \frac{1}{3} \quad P_{72}^{-2} = \frac{1}{3} \quad P_{71}^{-3} = \frac{1}{3}$$

$$P_{74}^{-1} = \frac{1}{3} \quad P_{73}^{-2} = \frac{1}{3} \quad P_{72}^{-3} = \frac{1}{3}$$

$$P_{75}^{-1} = \frac{1}{3} \quad P_{74}^{-2} = \frac{1}{3} \quad P_{73}^{-3} = \frac{1}{3}$$

Stan 8:

$$P_{84}^{-1} = \frac{1}{3} \quad P_{83}^{-2} = \frac{1}{3} \quad P_{82}^{-3} = \frac{1}{3}$$

$$P_{85}^{-1} = \frac{1}{3} \quad P_{84}^{-2} = \frac{1}{3} \quad P_{83}^{-3} = \frac{1}{3}$$

$$P_{86}^{-1} = \frac{1}{3} \quad P_{85}^{-2} = \frac{1}{3} \quad P_{84}^{-3} = \frac{1}{3}$$

Stan 9:

$$P_{95}^{-1} = \frac{1}{3} \quad P_{94}^{-2} = \frac{1}{3} \quad P_{93}^{-3} = \frac{1}{3}$$

$$P_{96}^{-1} = \frac{1}{3} \quad P_{95}^{-2} = \frac{1}{3} \quad P_{94}^{-3} = \frac{1}{3}$$

$$P_{97}^{-1} = \frac{1}{3} \quad P_{96}^{-2} = \frac{1}{3} \quad P_{95}^{-3} = \frac{1}{3}$$

Stan 10:

$$P_{10,6}^{-1} = \frac{1}{3} \quad P_{10,5}^{-2} = \frac{1}{3} \quad P_{10,4}^{-3} = \frac{1}{3}$$

$$P_{10,7}^{-1} = \frac{1}{3} \quad P_{10,6}^{-2} = \frac{1}{3} \quad P_{10,5}^{-3} = \frac{1}{3}$$

$$P_{10,8}^{-1} = \frac{1}{3} \quad P_{10,7}^{-2} = \frac{1}{3} \quad P_{10,6}^{-3} = \frac{1}{3}$$

## 2.1 Równania Bellmana

Zakładamy:  $\gamma = 1$ .

$$V^{k+1}(x) = \max_a \sum_y P_{xy}^a \cdot [R_{xy}^a + \gamma \cdot V^k(y)] \quad (2)$$

$$\begin{aligned}
V^{k+1}(0) &= \max\{1 \cdot (0 + \gamma \cdot V^k(0)), \quad 1 \cdot (0 + \gamma \cdot V^k(0)), \quad 1 \cdot (0 + \gamma \cdot V^k(0))\} \\
V^{k+1}(1) &= \max\{1 \cdot (0 + \gamma \cdot V^k(0)), \quad 1 \cdot (0 + \gamma \cdot V^k(0)), \quad 1 \cdot (0 + \gamma \cdot V^k(0))\} \\
V^{k+1}(2) &= \max\{1 \cdot (1 + \gamma \cdot V^k(0)), \quad 1 \cdot (0 + \gamma \cdot V^k(0)), \quad 1 \cdot (0 + \gamma \cdot V^k(0))\} \\
V^{k+1}(3) &= \max\{\frac{2}{3} \cdot (0 + \gamma \cdot V^k(0)) + \frac{1}{3} \cdot (0 + \gamma \cdot V^k(1)), \quad 1 \cdot (1 + \gamma \cdot V^k(0)), \quad 1 \cdot (0 + \gamma \cdot V^k(0))\} \\
V^{k+1}(4) &= \max\{\frac{1}{3} \cdot (0 + \gamma \cdot V^k(0)) + \frac{1}{3} \cdot (0 + \gamma \cdot V^k(1)) + \frac{1}{3} \cdot (0 + \gamma \cdot V^k(2)), \\
&\quad \frac{2}{3} \cdot (0 + \gamma \cdot V^k(0)) + \frac{1}{3} \cdot (0 + \gamma \cdot V^k(1)), \\
&\quad 1 \cdot (1 + \gamma \cdot V^k(0))\} \\
V^{k+1}(5) &= \max\{\frac{1}{3} \cdot (0 + \gamma \cdot V^k(1)) + \frac{1}{3} \cdot (0 + \gamma \cdot V^k(2)) + \frac{1}{3} \cdot (0 + \gamma \cdot V^k(3)), \\
&\quad \frac{1}{3} \cdot (0 + \gamma \cdot V^k(0)) + \frac{1}{3} \cdot (0 + \gamma \cdot V^k(1)) + \frac{1}{3} \cdot (0 + \gamma \cdot V^k(2)), \\
&\quad \frac{2}{3} \cdot (0 + \gamma \cdot V^k(0)) + \frac{1}{3} \cdot (0 + \gamma \cdot V^k(1))\} \\
V^{k+1}(6+i) &= \max\{\frac{1}{3} \cdot (0 + \gamma \cdot V^k(i+2)) + \frac{1}{3} \cdot (0 + \gamma \cdot V^k(i+3)) + \frac{1}{3} \cdot (0 + \gamma \cdot V^k(i+4)), \\
&\quad \frac{1}{3} \cdot (0 + \gamma \cdot V^k(i+1)) + \frac{1}{3} \cdot (0 + \gamma \cdot V^k(i+2)) + \frac{1}{3} \cdot (0 + \gamma \cdot V^k(i+3)), \\
&\quad \frac{1}{3} \cdot (0 + \gamma \cdot V^k(i)) + \frac{1}{3} \cdot (0 + \gamma \cdot V^k(i+1)) + \frac{1}{3} \cdot (0 + \gamma \cdot V^k(i+2))\}
\end{aligned}$$

Dla  $i \in \{0, \dots, N-6\}$ .

## 2.2 Wyznaczone wartości strategii

Wyznaczono funkcję wartości dla strategii:  $\pi_1(x) = -1$ ,  $\pi_2(x) = -2$ ,  $\pi_3(x) = -3$ . Na podstawie zebranych wyników, nie można określić, która ze strategii jest najlepsza.

$x$	$V^{\pi_1}(x)$	$V^{\pi_2}(x)$	$V^{\pi_3}(x)$
0	0.000000	0.000000	0.000000
1	0.000000	0.000000	0.000000
2	1.000000	0.000000	0.000000
3	0.000000	1.000000	0.000000
4	0.333333	0.000000	1.000000
5	0.333333	0.000000	0.000000
6	0.444444	0.333333	0.000000
7	0.222222	0.333333	0.000000
8	0.370370	0.333333	0.333333
9	0.333333	0.111111	0.333333
10	0.345679	0.222222	0.333333

**Tabela 1:** Funkcje wartości dla strategii

## 2.3 Optymalna strategia gry

Przy założeniu, że nasz przeciwnik gra w sposób losowy wyznaczono optymalną strategię gry.

$x$	$Q^{\pi^*}(x, -1)$	$Q^{\pi^*}(x, -2)$	$Q^{\pi^*}(x, -3)$
0	0.000000	0.000000	0.000000
1	0.000000	0.000000	0.000000
2	<u>1.000000</u>	0.000000	0.000000
3	0.000000	<u>1.000000</u>	0.000000
4	0.333333	0.000000	<u>1.000000</u>
5	<u>0.666667</u>	0.333333	0.000000
6	<u>1.000000</u>	0.666667	0.333333
7	0.888889	<u>1.000000</u>	0.666667
8	0.888889	0.888889	<u>1.000000</u>
9	0.888889	0.888889	0.888889
10	<u>1.000000</u>	0.888889	0.888889

**Tabela 2:** Funkcja wartości akcji

Optymalna funkcja wartości :  $V^*(x) = [1.0, 0.89, 1, 1, 1, 0.66, 1, 1, 1, 0, 0]$

Optymalna strategia wyznaczona jako zachłanna może być opisana sekwencją akcji:  
-1, \*, -3, -2, -1, -1, -3, -2, -1, \*, \*.

## 2.4 Kod źródłowy

Kod źródłowy zamieszczono również w formie tekstowej do zadania.

**Kod źródłowy 1:** Zadanie 3 - zapałki

Źródło: Opracowanie własne

```
1  """
2  Gra w zapałki
3
4  Mamy do dyspozycji pewną liczbę N zapałek.
5  W grze bierze udział dwóch graczy, którzy na przemian zabierają
6  1, 2 lub 3 zapałki.
7  Przegrywa ten gracz, który zabiera zapałki jako ostatni.
8  Przeciwnik będzie grał w sposób losowy.
9  """
10
11 import numpy as np
12
13 N = 10
14
15 def Q_next(V, gamma):
16     Q = np.zeros((N+1, 3))
```

```

17
18     # Stan 0
19     Q[0, 0] = 1 * (0 + gamma * V[0])
20     Q[0, 1] = 1 * (0 + gamma * V[0])
21     Q[0, 2] = 1 * (0 + gamma * V[0])
22
23     # Stan 1
24     Q[1, 0] = 1 * (0 + gamma * V[0])
25     Q[1, 1] = 1 * (0 + gamma * V[0])
26     Q[1, 2] = 1 * (0 + gamma * V[0])
27
28     # Stan 2
29     Q[2, 0] = 1 * (1 + gamma * V[0]) # ruch zwycięski
30     Q[2, 1] = 1 * (0 + gamma * V[0])
31     Q[2, 2] = 1 * (0 + gamma * V[0])
32
33     # Stan 3
34     Q[3, 0] = 2/3 * (0 + gamma * V[0]) + 1/3 * (0 + gamma * V[1])
35     Q[3, 1] = 1 * (1 + gamma * V[0]) # ruch zwycieski
36     Q[3, 2] = 1 * (0 + gamma * V[0])
37
38     # Stan 4
39     Q[4, 0] = 1/3 * (0 + gamma * V[0]) + 1/3 * (0 + gamma * V[1]) + 1/3 * (0 + gamma
    * V[2])
40     Q[4, 1] = 2/3 * (0 + gamma * V[0]) + 1/2 * (0 + gamma * V[1])
41     Q[4, 2] = 1 * (1 + gamma * V[0]) # ruch zwycieski
42
43     # Stan 5
44     Q[5, 0] = 1/3 * (0 + gamma * V[1]) + 1/3 * (0 + gamma * V[2]) + 1/3 * (0 + gamma
    * V[3])
45     Q[5, 1] = 1/3 * (0 + gamma * V[0]) + 1/3 * (0 + gamma * V[1]) + 1/3 * (0 + gamma
    * V[2])
46     Q[5, 2] = 2/3 * (0 + gamma * V[0]) + 1/3 * (0 + gamma * V[1])
47
48     # Stany od 6 do N
49     x = 6
50     for i in range(N-5):
51         Q[x + i, 0] = 1/3 * (0 + gamma * V[i+2]) + 1/3 * (0 + gamma * V[i+3]) + 1/3
    * (0 + gamma * V[i+4])
52         Q[x + i, 1] = 1/3 * (0 + gamma * V[i+1]) + 1/3 * (0 + gamma * V[i+2]) + 1/3
    * (0 + gamma * V[i+3])
53         Q[x + i, 2] = 1/3 * (0 + gamma * V[i]) + 1/3 * (0 + gamma * V[i+1]) + 1/3 *
    (0 + gamma * V[i+2])
54
55     return Q
56
57 def wartosc_strategii(id):
58     """
59     :param id: numer strategii gdzie 0 = -1, 1 = -2, 2 = -3
60     """
61     V_initial = np.zeros(N+1)

```

```

62     V = V_initial
63     V_prev = V
64
65     for i in range(100):
66         Q = Q_next(V, gamma=1)
67         V = Q[:, id]
68
69         if np.sum(V_prev - V) == 0:
70             break
71
72         V_prev = V
73
74     return Q
75
76 def optymalna_strategia():
77     V_initial = np.zeros(N+1)
78     V = V_initial
79     V_prev = V
80
81     for i in range(100):
82         Q = Q_next(V, gamma=1)
83         V = np.max(Q, axis=1)
84
85         if np.sum(V_prev - V) == 0:
86             break
87
88         V_prev = V
89
90     return Q
91
92
93 if __name__ == "__main__":
94     import pandas as pd
95
96     # Porównanie wartości strategii
97     print("V pi1")
98     Q1 = wartosc_strategii(0)
99     print(Q1[:, 0])
100    print("V pi2")
101    Q2 = wartosc_strategii(1)
102    print(Q2[:, 1])
103    print("V pi3")
104    Q3 = wartosc_strategii(2)
105    print(Q3[:, 2])
106
107    data = {"pi1": Q1[:, 0], "pi2": Q2[:, 1], "pi3": Q3[:, 2]}
108    df = pd.DataFrame.from_dict(data)
109    print(df.to_latex())
110
111
112    # Optymalna strategia

```

```

113     print("Q optymalnej strategii")
114     Q = optymalna_strategia()
115     df = pd.DataFrame(Q)
116     print(df.to_latex())

```

### 3 Zadanie 4

Przedmiotem rozważań w tym zadaniu była gra w ruletkę, gdzie obstawia się wielokrotności 1 zł na kolor z prawdopodobieństwem wygranej  $p = \frac{18}{37}$ .

Należało wyznaczyć optymalną strategię obstawiania, tak aby wygrać dokładnie 100 zł. Wykreślono optymalną funkcję wartości  $V^*(x)$  i optymalnej strategii  $\pi^*(x)$ .

#### 3.1 Założenia

Stan  $x$  określa kapitał gracza w pełnych zł:  $x \in \{0, 1, 2, 3, \dots, 100\}$ . Możliwe akcje do wykonania w stanie  $x$  to  $a \in \{1, \dots, \min(x, 100 - x)\}$ .

Oczekiwana wartość wzmocnienia:

$$R_{xy}^a = \begin{cases} 0 & \text{dla } y \neq 100 \\ 1 & \text{dla } y = 100 \end{cases} \quad (3)$$

Prawdopodobieństwa przejścia:

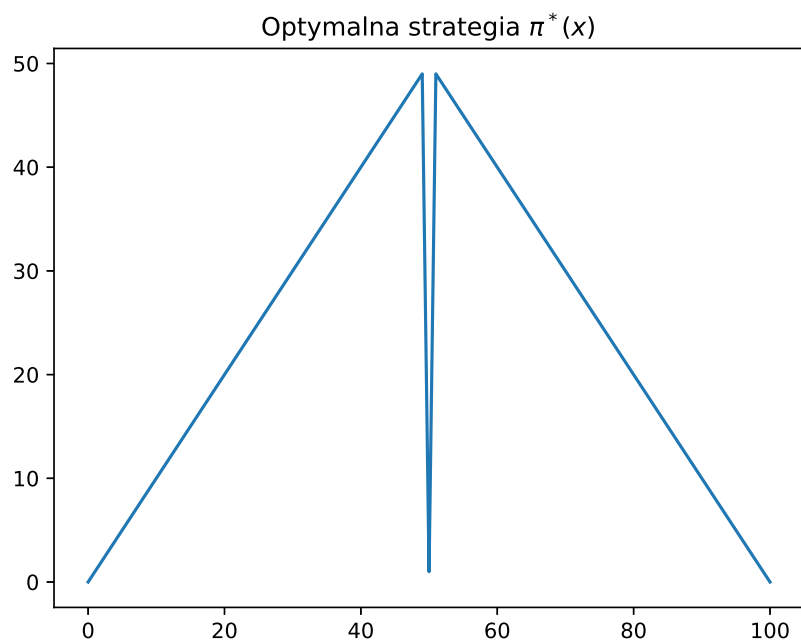
$$P_{xy}^a = \begin{cases} p & \text{dla } y = x + a \\ (1 - p) & \text{dla } y = x - a \end{cases} \quad (4)$$

#### 3.2 Wyniki

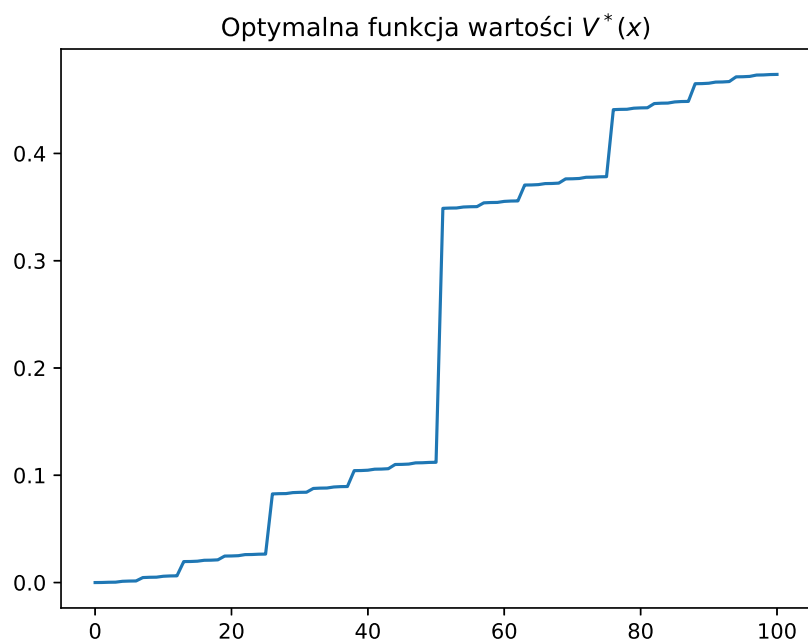
Wyznaczono wykres optymalnej funkcji wartości i optymalnej strategii korzystając z kodu z listingu 2.

Otrzymane wyniki są dosyć zaskakujące i trudno mi je uzasadnić. Optymalną strategią jest gra *va banque*, ryzykując wszystko. Jest to logiczne, ponieważ im dłużej gramy, tym częściej będziemy przegrywać. Nie rozumiem skąd wynika akcja 1 w stanie 50 zł. Być może gdzieś leży błąd w stworzonych równania Bellmana.

Optymalna funkcja wartości, przedstawiona na rysunku 2 przedstawia gwałtowne skoki w punktach takich jak 25, 50, 75. Tendencja rosnąca funkcji wartości jest zrozumiała, bo im więcej mamy pieniędzy (większy stan) tym bardziej prawdopodobne, że dojdziemy do kapitału 100 zł.



**Rysunek 1:** Optymalna strategia



**Rysunek 2:** Funkcja wartości

### 3.3 Kod źródłowy

**Kod źródłowy 2:** Zadanie 4 - ruletka

Źródło: Opracowanie własne



```

1  """
2  Gra w ruletkę.
3  Obstawiamy wielokrotność 1 zł na kolor.
4  Porawdopodobieństwo wygranej  $p = 18/37$ 
5
6  Wyznaczyć optymalną strategię obstawiania, tak aby wygrać dokładnie 100 zł.
7  Sporządzić wykres optymalnej funkcji wartości  $V^*(x)$  i optymalnej strategii  $\pi^*(x)$ 
8
9  Stan  $x$  określa kapitał gracza  $x$  in  $\{0, 1, 2, 3, \dots, 100\}$ 
10 Możliwe akcje  $a$  w stanie  $x$ :  $a$  in  $\{1, \dots, \min(x, 100-x)\}$ 
11     Czyli dla maksymalna akcja = 50 (bo dla  $x = 50$  możemy postawić maksymalnie 50).
12     Co dla stanu 0? Raczej nie ma żadnej możliwej akcji.
13 Oczekiwana wartość wzmocnienia:  $R_{\{xy\}} = 1$  dla  $y = 100$ , w innych przypadkach  $R=0$ 
14 Prawodopodobieństwa przejścia  $p$  dla  $y = x+a$  (wygrana), w innym przypadku  $(1-p)$ 
15 """
16
17
18 import numpy as np
19 import matplotlib.pyplot as plt
20
21 p = 18 / 37
22
23
24 def Q_next(V, gamma):
25     Q = np.zeros((101, 50)) # każdy maksymalnie będzie 50 stanów (dla  $x = 50$ 
    najwiecej możliwych akcji do wykonania)
26
27     for x in range(101): # dla wszystkich stanów
28         for a in range(50): # dla wszystkich akcji
29             if a > min(x, 100 - x): # niemożliwe akcje
30                 Q[x, a] = 0
31             else:
32                 loss = (1 - p) * (0 + gamma * V[x - a])
33                 R = 1 if x + a == 100 else 0
34                 win = p * (R + gamma * V[x + a])
35                 Q[x, a] = p * win + (1 - p) * loss
36
37     return Q
38
39
40 def optymalna_strategia():
41     V_initial = np.zeros(101) # dla każdego stanu  $\{0, 1, 2, \dots, 100\}$ 
42     V = V_initial
43     V_prev = V
44
45     while True:
46         Q = Q_next(V, gamma=1)
47         V = np.max(Q, axis=1)
48
49         if np.sum(V_prev - V) == 0:
50             break

```

```

51
52     V_prev = V
53
54     return Q
55
56
57 if __name__ == "__main__":
58     print("Q optymalnej strategii")
59     Q = optymalna_strategia()
60     print(Q[50, :])
61
62     # TODO: coś mi nie pasuje, bo:
63     # występują skoki
64     # dla stanu 50 zł powinno się postawić złotówkę, w innych przypadkach trzeba
    grać va bank
65     # dla stanu 50 zł wartość akcji jest wszędzie prawie identyczna
66
67     V = np.max(Q, axis=1)
68     plt.plot(V)
69     plt.title("Optymalna funkcja wartości  $V^*(x)$ ")
70     plt.show()
71     pi = np.argmax(Q, axis=1)
72     plt.plot(pi)
73     plt.title("Optymalna strategia  $\pi^*(x)$ ")
74     plt.show()
75     print(Q)

```