



Федеральное государственное бюджетное образовательное учреждение
высшего профессионального образования
«Московский государственный технический университет имени Н.Э.
Баумана»
(МГТУ им. Н.Э. Баумана)

**Лабораторная работа №4
по курсу «Методы машинного обучения»**

Выполнил
студент группы ИУ5-22М
XXXX

Москва, 2023

1. Задание

1. На основе рассмотренного на лекции примера реализуйте алгоритм Policy Iteration для любой среды обучения с подкреплением (кроме рассмотренной на лекции среды Toy Text / Frozen Lake) из библиотеки Gym (или аналогичной библиотеки).

2. Текст программы

```
1 #!/usr/bin/env python
2
3 import gymnasium as gym
4 import numpy as np
5 from pprint import pprint
6
7 class PolicyIterationAgent:
8     '''
9     Класс, эмулирующий работу агента
10    '''
11    def __init__(self, env):
12        self.env = env
13        # Пространство состояний
14        self.observation_dim = 500
15        # Массив действий в соответствии с документацией
16        # https://gymnasium.farama.org/environments/toy_text/taxi/
17        self.actions_variants = np.array([0,1,2,3,4,5])
18        # 0: Move south (down)
19        # 1: Move north (up)
20        # 2: Move east (right)
21        # 3: Move west (left)
22        # 4: Pickup passenger
23        # 5: Drop off passenger
24        # Задание стратегии (политики)
25        # Карта 5x5 и 6 возможных действия
26        self.policy_probs = np.full((self.observation_dim, len(self.actions_variants)), 0.25)
27        # Начальные значения для v(s)
28        self.state_values = np.zeros(shape=(self.observation_dim))
29        # Начальные значения параметров
30        self.maxNumberOfIterations = 1000
31        self.theta=1e-6
32        self.gamma=0.99
33
34    def print_policy(self):
35        '''
36        Вывод матриц стратегии
37        '''
38        print('Стратегия:')
39        pprint(self.policy_probs)
40
41    def policy_evaluation(self):
42        '''
43        Оценивание стратегии
44        '''
45        # Предыдущее значение функции ценности
46        valueFunctionVector = self.state_values
47        for iterations in range(self.maxNumberOfIterations):
48            # Новое значение функции ценности
49            valueFunctionVectorNextIteration=np.zeros(shape=(self.observation_dim))
```

```

50     # Цикл по состояниям
51     for state in range(self.observation_dim):
52         # Вероятности действий
53         action_probabilities = self.policy_probs[state]
54         # Цикл по действиям
55         outerSum=0
56         for action, prob in enumerate(action_probabilities):
57             innerSum=0
58             # Цикл по вероятностям действий
59             for probability, next_state, reward, isTerminalState in
self.env.P[state][action]:
60                 innerSum=innerSum+probability*(reward+self.gamma*self.state_values[next_
t_state])
61             outerSum=outerSum+self.policy_probs[state][action]*innerSum
62             valueFunctionVectorNextIteration[state]=outerSum
63             if (np.max(np.abs(valueFunctionVectorNextIteration-valueFunctionVector))<self.theta_
):
64                 # Проверка сходимости алгоритма
65                 valueFunctionVector=valueFunctionVectorNextIteration
66                 break
67             valueFunctionVector=valueFunctionVectorNextIteration
68     return valueFunctionVector
69
70
71 def policy_improvement(self):
72     '''
73     Улучшение стратегии
74     '''
75     qvaluesMatrix=np.zeros((self.observation_dim, len(self.actions_variants)))
76     improvedPolicy=np.zeros((self.observation_dim, len(self.actions_variants)))
77     # Цикл по состояниям
78     for state in range(self.observation_dim):
79         for action in range(len(self.actions_variants)):
80             for probability, next_state, reward, isTerminalState in
self.env.P[state][action]:
81                 qvaluesMatrix[state,action]=qvaluesMatrix[state,action]+probability*(rewar
d+self.gamma*self.state_values[next_state])
82
83         # Находим лучшие индексы
84         bestActionIndex=np.where(qvaluesMatrix[state,:]==np.max(qvaluesMatrix[state,:]))
85         # Обновление стратегии
86         improvedPolicy[state,bestActionIndex]=1/np.size(bestActionIndex)
87     return improvedPolicy
88
89
90 def policy_iteration(self, cnt):
91     '''
92     Основная реализация алгоритма
93     '''
94     policy_stable = False
95     for i in range(1, cnt+1):
96         self.state_values = self.policy_evaluation()
97         self.policy_probs = self.policy_improvement()
98         print(f'{i} шагов.')
99
100 def play_agent(agent):
101     env2 = gym.make('Taxi-v3', render_mode='human')
102     state = env2.reset()[0]
103     done = False
104     while not done:
105         p = agent.policy_probs[state]

```

```

106         if isinstance(p, np.ndarray):
107             action = np.random.choice(len(agent.actions_variants), p=p)
108         else:
109             action = p
110         next_state, reward, terminated, truncated, _ = env2.step(action)
111         env2.render()
112         state = next_state
113         if terminated or truncated:
114             done = True
115
116 def main():
117     # Создание среды
118     env = gym.make('Taxi-v3')
119     env.reset()
120     # Обучение агента
121     agent = PolicyIterationAgent(env)
122     agent.policy_iteration(1000)
123     agent.print_policy()
124     # Проигрывание сцены для обученного агента
125     play_agent(agent)
126
127 if __name__ == '__main__':
128     main()

```

3. Экранные формы с примерами выполнения программы

```

1 $ ./4.py
2 1000 шагов.
3 Стратегия:
4 array([[0. , 0. , 0. , 0. , 1. , 0. ],
5        [0. , 0. , 0. , 0. , 1. , 0. ],
6        [0. , 0. , 0. , 0. , 1. , 0. ],
7        ...,
8        [0. , 1. , 0. , 0. , 0. , 0. ],
9        [0. , 0.5, 0. , 0.5, 0. , 0. ],
10       [0. , 0. , 0. , 1. , 0. , 0. ]])

```

