

МОСКОВСКИЙ ГОСУДАРСТВЕННЫЙ ТЕХНИЧЕСКИЙ
УНИВЕРСИТЕТ
им. Н.Э. Баумана

Факультет «Информатика и системы управления»
Кафедра «Систем обработки информации и управления»

ОТЧЕТ

Лабораторная работа № 4
по дисциплине «Методы машинного обучения»

Тема: «Реализация алгоритма Policy Iteration.»

ИСПОЛНИТЕЛЬ: Кожуро Б.Е. ФИО
группа ИУ5-21М
подпись
" " 2024 г.

ПРЕПОДАВАТЕЛЬ: Гапанюк Ю.Е. ФИО
подпись
" " 2024 г.

Москва - 2024

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. Сформировать отчет и разместить его в своем репозитории на github.

Выполнение

Для реализации была выбрана среда Taxi-v3 из библиотеки Gym.

По документации: 500 состояний – 5*5 карта, 4 возможных локации точки выхода, 5 состояний пассажира (4 выхода и в такси).

6 действий – 4 движения и взять/высадить пассажира.

Из 500 состояний в рамках 1 итерации достижимо 400 – исключаются состояния где пассажир там же, где и здание.

Код программы:

```
import gym
import numpy as np
import matplotlib.pyplot as plt
from pprint import pprint
import pandas as pd
from gym.envs.toy_text.taxi import TaxiEnv

def print_full(x):
    pd.set_option('display.max_rows', len(x))
    print(x)
    pd.reset_option('display.max_rows')

class PolicyIterationAgent:
    """
    Класс, эмулирующий работу агента
    """
    def __init__(self, env):
        self.env = env
        # Пространство состояний
        self.observation_dim = 500
        # Массив действий в соответствии с документацией
        self.actions_variants = np.array([0,1,2,3,4,5])
        # Задание стратегии (политики)
        self.policy_probs = np.full((self.observation_dim,
len(self.actions_variants)), 0.16666666)
        # Начальные значения для v(s)
        self.state_values = np.zeros(shape=(self.observation_dim))
        # Начальные значения параметров
```

```

self.maxNumberOfIterations = 1000
self.theta=1e-6
self.gamma=0.99

def print_policy(self):
    '''
    Вывод матриц стратегии
    '''

    if self.policy_probs[0][0] != 0.16666666:
        #np.set_printoptions(threshold=np.inf)
        x = TaxiEnv()
        pos = {0:'R', 1:'G',2:'Y', 3:'B', 4:'T'}
        print(''
+-----+
|R: | : :G|
| : | : : |
| : : : : |
| | : | : |
|Y| : |B: |
+-----+
        ''')
        print('состояние: x,y,пассажир,назначение')
        print('Стратегия:')
        for i in range(len(self.policy_probs)):
            t_x,t_y,passeng,dest = x.decode(i)
            print((t_x,t_y,pos[passeng],pos[dest]), self.policy_probs[i])
            #np.set_printoptions(threshold=False)
    else:
        print('Стратегия:')
        pprint(self.policy_probs)

def policy_evaluation(self):
    '''
    Оценивание стратегии
    '''

    # Предыдущее значение функции ценности
    valueFunctionVector = self.state_values
    for iterations in range(self.maxNumberOfIterations):
        # Новое значение функции ценности
        valueFunctionVectorNextIteration=np.zeros(shape=(self.observation_dim
))

        # Цикл по состояниям
        for state in range(self.observation_dim):
            # Вероятности действий
            action_probabilities = self.policy_probs[state]
            # Цикл по действиям
            outerSum=0
            for action, prob in enumerate(action_probabilities):
                innerSum=0

```

```

        # Цикл по вероятностям действий
        for probability, next_state, reward, isTerminalState in
self.env.P[state][action]:
            innerSum=innerSum+probability*(reward+self.gamma*self.state_values[next_state])
            outerSum=outerSum+self.policy_probs[state][action]*innerSum
            valueFunctionVectorNextIteration[state]=outerSum
            if(np.max(np.abs(valueFunctionVectorNextIteration-valueFunctionVector))<self.theta):
                # Проверка сходимости алгоритма
                valueFunctionVector=valueFunctionVectorNextIteration
                break
            valueFunctionVector=valueFunctionVectorNextIteration
        return valueFunctionVector

def policy_improvement(self):
    """
    Улучшение стратегии
    """
    qvaluesMatrix=np.zeros((self.observation_dim,
len(self.actions_variants)))
    improvedPolicy=np.zeros((self.observation_dim,
len(self.actions_variants)))
    # Цикл по состояниям
    for state in range(self.observation_dim):
        for action in range(len(self.actions_variants)):
            for probability, next_state, reward, isTerminalState in
self.env.P[state][action]:
                qvaluesMatrix[state,action]=qvaluesMatrix[state,action]+probability*(reward+self.gamma*self.state_values[next_state])

            # Находим лучшие индексы
            bestActionIndex=np.where(qvaluesMatrix[state,:]==np.max(qvaluesMatrix[state,:]))
            # Обновление стратегии
            improvedPolicy[state,bestActionIndex]=1/np.size(bestActionIndex)
    return improvedPolicy

def policy_iteration(self, cnt):
    """
    Основная реализация алгоритма
    """
    policy_stable = False
    for i in range(1, cnt+1):
        self.state_values = self.policy_evaluation()
        self.policy_probs = self.policy_improvement()
        print(f'Алгоритм выполнен за {i} шагов.')

def play_agent(agent):

```

```

env2 = gym.make('Taxi-v3', render_mode='human')
state = env2.reset()[0]
done = False
while not done:
    p = agent.policy_probs[state]
    if isinstance(p, np.ndarray):
        action = np.random.choice(len(agent.actions_variants), p=p)
    else:
        action = p
    next_state, reward, terminated, truncated, _ = env2.step(action)
    env2.render()
    state = next_state
    if terminated or truncated:
        done = True

def main():
    # Создание среды
    env = gym.make('Taxi-v3')
    env.reset()
    # Обучение агента
    agent = PolicyIterationAgent(env)
    agent.print_policy()
    agent.policy_iteration(1000)
    agent.print_policy()
    # Проигрывание сцены для обученного агента
    play_agent(agent)

if __name__ == '__main__':
    main()

```

Вывод программы модифицирован для кодировки состояний: см. рис. 2.

Результаты выполнения:

```

Стратегия:
array([[0.16666666, 0.16666666, 0.16666666, 0.16666666, 0.16666666,
        0.16666666],
       [0.16666666, 0.16666666, 0.16666666, 0.16666666, 0.16666666,
        0.16666666],
       [0.16666666, 0.16666666, 0.16666666, 0.16666666, 0.16666666,
        0.16666666],
       ...,
       [0.16666666, 0.16666666, 0.16666666, 0.16666666, 0.16666666,

```

Рис. 1 – начальная стратегия.

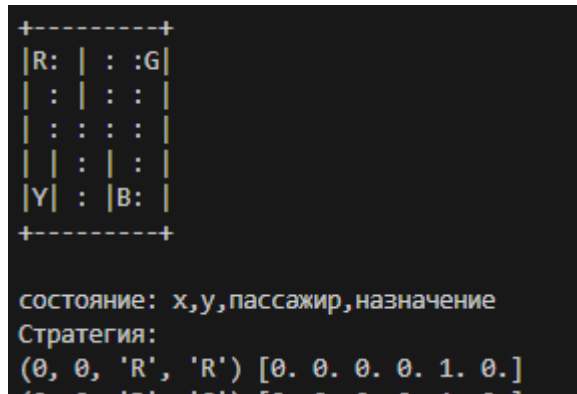


Рис. 2 – Модифицированный вывод.

Фрагмент вывода стратегии:

```

(0, 0, 'R', 'R') [0. 0. 0. 0. 1. 0.]
(0, 0, 'R', 'G') [0. 0. 0. 0. 1. 0.]
(0, 0, 'R', 'Y') [0. 0. 0. 0. 1. 0.]
(0, 0, 'R', 'B') [0. 0. 0. 0. 1. 0.]
(0, 0, 'G', 'R') [0.5 0. 0.5 0. 0. 0. ]
(0, 0, 'G', 'G') [0.5 0. 0.5 0. 0. 0. ]
(0, 0, 'G', 'Y') [0.5 0. 0.5 0. 0. 0. ]
(0, 0, 'G', 'B') [0.5 0. 0.5 0. 0. 0. ]
(0, 0, 'Y', 'R') [1. 0. 0. 0. 0. 0.]
(0, 0, 'Y', 'G') [1. 0. 0. 0. 0. 0.]
(0, 0, 'Y', 'Y') [1. 0. 0. 0. 0. 0.]
(0, 0, 'Y', 'B') [1. 0. 0. 0. 0. 0.]
(0, 0, 'B', 'R') [0.5 0. 0.5 0. 0. 0. ]
(0, 0, 'B', 'G') [0.5 0. 0.5 0. 0. 0. ]
(0, 0, 'B', 'Y') [0.5 0. 0.5 0. 0. 0. ]
(0, 0, 'B', 'B') [0.5 0. 0.5 0. 0. 0. ]
(0, 0, 'T', 'R') [0. 0. 0. 0. 0. 1.]
(0, 0, 'T', 'G') [0.5 0. 0.5 0. 0. 0. ]
(0, 0, 'T', 'Y') [1. 0. 0. 0. 0. 0.]
(0, 0, 'T', 'B') [0.5 0. 0.5 0. 0. 0. ]
(0, 1, 'R', 'R') [0. 0. 0. 1. 0. 0.]
(0, 1, 'R', 'G') [0. 0. 0. 1. 0. 0.]
(0, 1, 'R', 'Y') [0. 0. 0. 1. 0. 0.]
(0, 1, 'R', 'B') [0. 0. 0. 1. 0. 0.]
(0, 1, 'G', 'R') [1. 0. 0. 0. 0. 0.]
. . .
(4, 4, 'R', 'R') [0. 0.5 0. 0.5 0. 0. ]
(4, 4, 'R', 'G') [0. 0.5 0. 0.5 0. 0. ]
(4, 4, 'R', 'Y') [0. 0.5 0. 0.5 0. 0. ]
(4, 4, 'R', 'B') [0. 0.5 0. 0.5 0. 0. ]
(4, 4, 'G', 'R') [0. 1. 0. 0. 0. 0.]
(4, 4, 'G', 'G') [0. 1. 0. 0. 0. 0.]
(4, 4, 'G', 'Y') [0. 1. 0. 0. 0. 0.]
(4, 4, 'G', 'B') [0. 1. 0. 0. 0. 0.]
(4, 4, 'Y', 'R') [0. 0.5 0. 0.5 0. 0. ]
(4, 4, 'Y', 'G') [0. 0.5 0. 0.5 0. 0. ]
(4, 4, 'Y', 'Y') [0. 0.5 0. 0.5 0. 0. ]
(4, 4, 'Y', 'B') [0. 0.5 0. 0.5 0. 0. ]
(4, 4, 'B', 'R') [0. 0. 0. 1. 0. 0.]
(4, 4, 'B', 'G') [0. 0. 0. 1. 0. 0.]
(4, 4, 'B', 'Y') [0. 0. 0. 1. 0. 0.]
(4, 4, 'B', 'B') [0. 0. 0. 1. 0. 0.]
(4, 4, 'T', 'R') [0. 0.5 0. 0.5 0. 0. ]
(4, 4, 'T', 'G') [0. 1. 0. 0. 0. 0.]
(4, 4, 'T', 'Y') [0. 0.5 0. 0.5 0. 0. ]
(4, 4, 'T', 'B') [0. 0. 0. 1. 0. 0.]

```

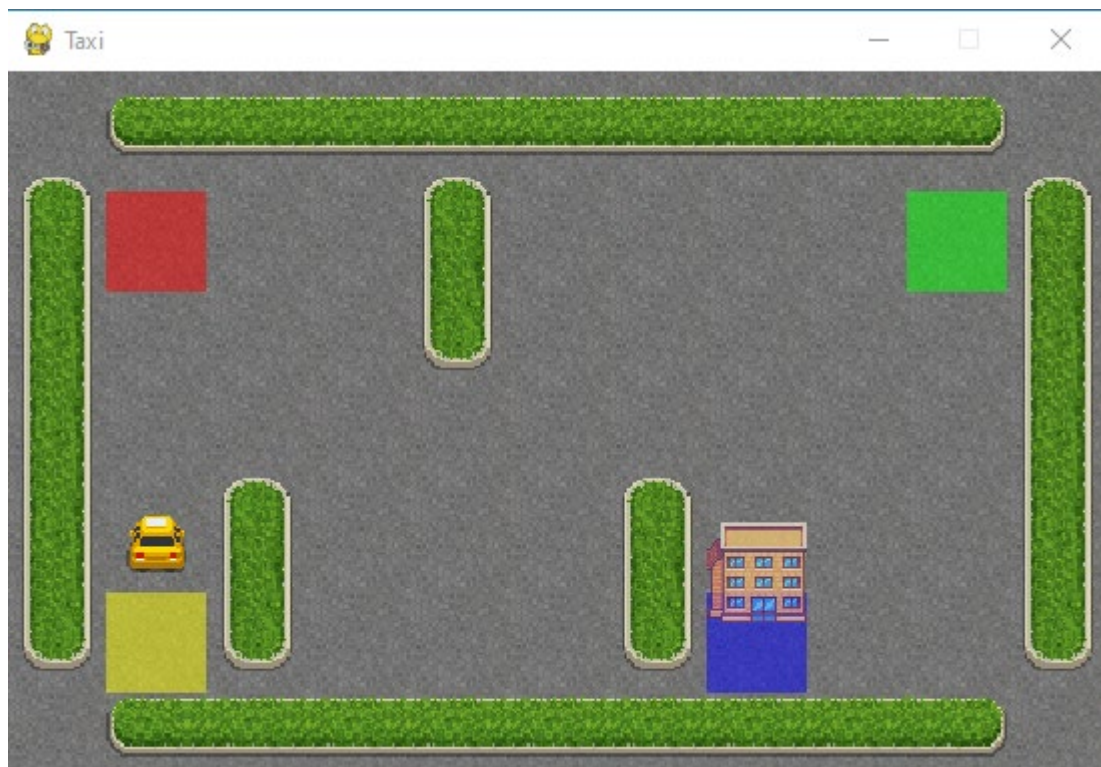


Рис. 3 – работа программы.

Вывод:

Методика Policy Iteration позволяет, имея матрицу состояний и вероятностей действий, итеративно улучшать стратегию переходов между состояниями. В данной ЛР улучшение достигается за счёт штрафа за лишние переходы и штрафов за взятие и высадку пассажира вне ожидаемой зоны.

Таким образом, все переходы будут равнозначны до момента нахождения пассажира, и его случайной доставки.

Оптимизация стратегии начинается с +20 для состояний, в которых пассажир в такси, и он на нужной точке высадки (результат действия $v = -10 + 20 = 10$), после этого соседние состояния, имеющие возможность попасть в это, будут направлены в него. Технически решение достижимо за 12 итераций для 1 точки высадки.