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:	<u> Савельев А.А.</u>
группа ИУ5-24М	ФИС
	""2024 г.
ПРЕПОДАВАТЕЛЬ:	<u>Гапанюк Ю.Е.</u>
	подпись
	""2024 г.

Москва - 2024

Лабораторная работа N°4 Реализация алгоритма Policy Iteration.

Цель лабораторной работы: ознакомление с базовыми методами обучения с подкреплением.

Задание:

• 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).

```
import gym
import numpy as np
import time
import matplotlib.pyplot as plt
from pprint import pprint
def main():
    state, action = 0, 0
    env = gym.make("CliffWalking-v0")
    print('Пространство состояний:')
    pprint(env.observation space)
    print()
    print('Пространство действий:')
    pprint(env.action space)
    print()
    print('Диапазон наград:')
    pprint(env.reward range)
    print()
    print('Вероятности для 0 состояния и 0 действия:')
    pprint(env.P[state][action])
    print('Вероятности для 0 состояния:')
    pprint(env.P[state])
if __name__ == '__main__':
    main()
Пространство состояний:
Discrete(48)
Пространство действий:
Discrete(4)
Диапазон наград:
(-inf, inf)
```

```
Вероятности для 0 состояния и 0 действия:
[(1.0, 0, -1, False)]
Вероятности для 0 состояния:
{0: [(1.0, 0, -1, False)],
1: [(1.0, 1, -1, False)],
2: [(1.0, 12, -1, False)],
3: [(1.0, 0, -1, False)]}
class PolicyIterationAgent:
    Класс, эмулирующий работу агента
    def __init__(self, env):
        self.env = env
        # Пространство состояний
        self.observation dim = 48
        # Массив действий в соответствии с документацией
https://www.gymlibrary.dev/environments/toy_text/frozen_lake/
        self.actions variants = np.array([0, 1, 2, 3])
        # Задание стратегии (политики)
        # Карта 4х4 и 6 возможных действий
        self.policy probs = np.full((self.observation dim,
len(self.actions variants)), 0.25)
        # 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
        print('Стратегия:')
        pprint(self.policy_probs)
    def policy evaluation(self):
        Оценивание стратегии
        1.1.1
        # 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):</pre>
                # Проверка сходимости алгоритма
                valueFunctionVector = valueFunctionVectorNextIteration
            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 выполнился за {i} шагов.')
def play agent(agent):
    env2 = gym.make('CliffWalking-v0', 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('CliffWalking-v0')
    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()
Стратегия:
array([[0.25, 0.25, 0.25, 0.25],
```

```
[0.25, 0.25, 0.25, 0.25],
       [0.25, 0.25, 0.25, 0.25],
       [0.25, 0.25, 0.25, 0.25],
       [0.25, 0.25, 0.25, 0.25],
       [0.25, 0.25, 0.25, 0.25],
       [0.25, 0.25, 0.25, 0.25],
       [0.25, 0.25, 0.25, 0.25],
       [0.25, 0.25, 0.25, 0.25],
       [0.25, 0.25, 0.25, 0.25],
       [0.25, 0.25, 0.25, 0.25],
       [0.25, 0.25, 0.25, 0.25],
       [0.25, 0.25, 0.25, 0.25],
       [0.25, 0.25, 0.25, 0.25],
       [0.25, 0.25, 0.25, 0.25],
       [0.25, 0.25, 0.25, 0.25],
       [0.25, 0.25, 0.25, 0.25],
       [0.25, 0.25, 0.25, 0.25],
       [0.25, 0.25, 0.25, 0.25],
       [0.25, 0.25, 0.25, 0.25],
       [0.25, 0.25, 0.25, 0.25],
       [0.25, 0.25, 0.25, 0.25],
       [0.25, 0.25, 0.25, 0.25],
       [0.25, 0.25, 0.25, 0.25],
       [0.25, 0.25, 0.25, 0.25],
       [0.25, 0.25, 0.25, 0.25],
       [0.25, 0.25, 0.25, 0.25],
       [0.25, 0.25, 0.25, 0.25],
       [0.25, 0.25, 0.25, 0.25],
       [0.25, 0.25, 0.25, 0.25],
       [0.25, 0.25, 0.25, 0.25],
       [0.25, 0.25, 0.25, 0.25],
       [0.25, 0.25, 0.25, 0.25],
       [0.25, 0.25, 0.25, 0.25],
       [0.25, 0.25, 0.25, 0.25],
       [0.25, 0.25, 0.25, 0.25],
       [0.25, 0.25, 0.25, 0.25],
       [0.25, 0.25, 0.25, 0.25],
       [0.25, 0.25, 0.25, 0.25],
       [0.25, 0.25, 0.25, 0.25],
       [0.25, 0.25, 0.25, 0.25],
       [0.25, 0.25, 0.25, 0.25],
       [0.25, 0.25, 0.25, 0.25],
       [0.25, 0.25, 0.25, 0.25],
       [0.25, 0.25, 0.25, 0.25],
       [0.25, 0.25, 0.25, 0.25],
       [0.25, 0.25, 0.25, 0.25],
       [0.25, 0.25, 0.25, 0.25]])
Алгоритм выполнился за 1000 шагов.
Стратегия:
```

```
array([[0. , 0.5 , 0.5 , 0. [0.333333333, 0.333333333, 0.333333333, 0.
     [0. , 0. , 1. , 0. ]
[0. , 0. , 1. , 0. ]
[0. , 0. , 1. , 0. ]
[0. , 0. , 1. , 0. ]
    [0.
     , 0.33333333, 0.33333333, 0.33333333],
    [0.33333333, 0.33333333, 0.33333333, 0.
```

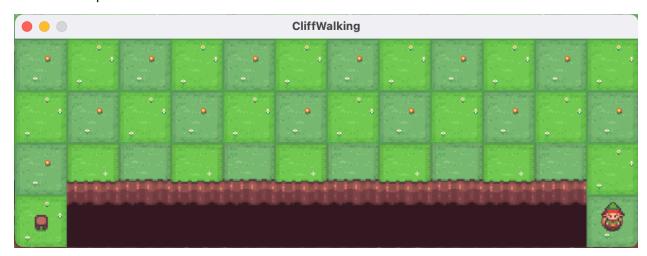
При выполнении кода в текущей ячейке или предыдущей ячейке ядро аварийно завершило работу.

Проверьте код в ячейках, чтобы определить возможную причину сбоя.

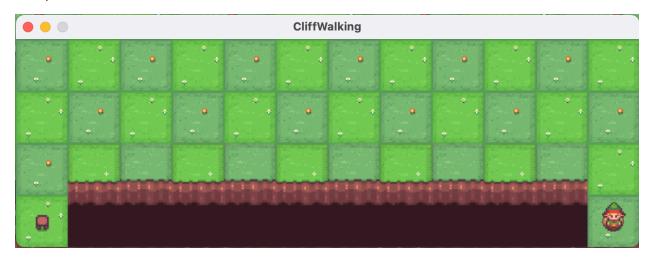
Щелкните здесь, чтобы получить дополнительные сведения.

Подробнее см. в журнале Jupyter.

В начале алгоритма



В конце



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 тортика.