



**Министерство науки и высшего образования Российской Федерации
Федеральное государственное бюджетное образовательное учреждение
высшего образования
«Московский государственный технический университет
имени Н.Э. Баумана
(национальный исследовательский университет)»
(МГТУ им. Н.Э. Баумана)**

**Факультет «Информатика и системы управления»
Кафедра ИУ5 «Системы обработки информации и управления»**

Курс «Методы машинного обучения»

Отчет по домашнему заданию

Выполнила:
студент группы ИУ5-22М
Павловская А.А.
23.05.2023

Проверил:
преподаватель каф. ИУ5
Гапанюк Ю.Е.

Москва, 2023 г.

Цель работы: ознакомление с методом обучения с подкреплением на основе временных различий Double SARSA.

Задание:

Реализовать алгоритм Double SARSA для среды обучения с подкрепления CliffWalking из библиотеки Gym

Ход работы

В среде CliffWalking агент может находиться в 48 состояниях и осуществлять 4 действия. Одно состояние начальное, одно конечное, а также есть «обрыв» из 10 клеток внизу карты (рис.1).

1	2	3	4	5	6	7	8	9	10	11	12
13	14	15	16	17	18	19	20	21	22	23	24
25	26	27	28	29	30	31	32	33	34	35	36
37	38	39	40	41	42	43	44	45	46	47	48

Рис.1. Карта состояний

Действия: 1 – вверх, 2 – вправо, 3 – вниз, 4 – влево.

Вознаграждение равно -1 для всех переходов, кроме ведущих в «обрыв». За вход в эту область начисляется вознаграждение -100 , после чего агент мгновенно переносится в начальную позицию [2].

В алгоритме Double SARSA, как и в Double Q-learning, используется эпсилон-жадная стратегия, определяющая действие на основе среднего по сумме значений двух таблиц Q^1 и Q^2 .

Обновление происходит по правилу [1]:

$$Q_{t+1}^1(s, a) \leftarrow Q_t^1(s, a) + \alpha[r + \gamma Q_t^2(s', a') - Q_t^1(s, a)]$$

Алгоритм Double SARSA отличается от SARSA наличием дополнительной таблицы, а от Double Q-learning тем, что в конце Q^1 и Q^2 не обновляются с вероятностью 0.5, а обновляется лишь одна таблица, а затем с вероятностью 0.5 они меняются местами.

Алгоритм 1. Double SARSA

1. Инициализировать $Q^1(s, a)$ и $Q^2(s, a)$ произвольным образом
2. **Повторять** для каждого эпизода:
 3. Инициализировать s
 4. Выбрать a в состоянии s , следуя произвольной стратегии
 5. **Повторять** для каждого шага эпизода:
 6. Предпринять действие, наблюдать r, s'
 7. Выбрать a' из состояния s' , следуя ϵ -жадной стратегии относительно среднего $Q^1 + Q^2$
 8. $Q^1(s, a) \leftarrow Q^1(s, a) + \alpha[r + \gamma Q^2(s', a') - Q^1(s, a)]$

9. $s \leftarrow s'; a \leftarrow a';$
10. С вероятностью 0.5:
11. Поменять местами Q^1 и Q^2

Текст программы

```
import numpy as np
import matplotlib.pyplot as plt
import gym
from tqdm import tqdm

# ***** БАЗОВЫЙ АГЕНТ
# *****

class BasicAgent:
    """
    Базовый агент, от которого наследуются стратегии обучения
    """

    # Наименование алгоритма
    ALGO_NAME = '---'

    def __init__(self, env, eps=0.1):
        # Среда
        self.env = env
        # Размерности Q-матрицы
        self.nA = env.action_space.n
        self.nS = env.observation_space.n
        #и сама матрица
        self.Q = np.zeros((self.nS, self.nA))
        # Значения коэффициентов
        # Порог выбора случайного действия
        self.eps=eps
        # Награды по эпизодам
        self.episodes_reward = []

    def print_q(self):
        print('Вывод Q-матрицы для алгоритма ', self.ALGO_NAME)
        print(self.Q)

    def get_state(self, state):
        """
        Возвращает правильное начальное состояние
        """
        if type(state) is tuple:
            # Если состояние вернулось с виде кортежа, то вернуть только номер
            # состояния
            return state[0]
```

```

        else:
            return state

def greedy(self, state):
    """
    <<Жадное>> текущее действие
    Возвращает действие, соответствующее максимальному Q-значению
    для состояния state
    """
    return np.argmax(self.Q[state])

def make_action(self, state):
    """
    Выбор действия агентом
    """
    if np.random.uniform(0,1) < self.eps:

        # Если вероятность меньше eps
        # то выбирается случайное действие
        return self.env.action_space.sample()
    else:
        # иначе действие, соответствующее максимальному Q-значению
        return self.greedy(state)

def draw_episodes_reward(self):
    # Построение графика наград по эпизодам
    fig, ax = plt.subplots(figsize = (15,10))
    y = self.episodes_reward
    x = list(range(1, len(y)+1))
    plt.plot(x, y, '-', linewidth=1, color='green')
    plt.title('Награды по эпизодам')
    plt.xlabel('Номер эпизода')
    plt.ylabel('Награда')
    plt.show()

def learn():
    """
    Реализация алгоритма обучения
    """
    pass

# ***** Double SARSA
# *****

class DoubleSARSA_Agent(BasicAgent):
    """
    Реализация алгоритма Double SARSA

```

```

...
# Наименование алгоритма
ALGO_NAME = 'Double SARSA'

def __init__(self, env, eps=0.015, lr=0.1, gamma=0.97, num_episodes=20000):
#self, env, eps=0.4, lr=0.1, gamma=0.98, num_episodes=20000
    # Вызов конструктора верхнего уровня
    super().__init__(env, eps)
    # Вторая матрица
    self.Q2 = np.zeros((self.nS, self.nA))
    # Learning rate
    self.lr=lr
    # Коэффициент дисконтирования
    self.gamma = gamma
    # Количество эпизодов
    self.num_episodes=num_episodes
    # Постепенное уменьшение eps
    self.eps_decay=0.00005
    self.eps_threshold=0.01

def greedy(self, state):
    ...

    <<Жадное>> текущее действие
    Возвращает действие, соответствующее максимальному Q-значению
    для состояния state
    ...

    temp_q = self.Q[state] + self.Q2[state]
    return np.argmax(temp_q)

def print_q(self):
    print('Вывод Q-матриц для алгоритма ', self.ALGO_NAME)
    print('Q1')
    print(self.Q)
    print('Q2')
    print(self.Q2)

def learn(self):
    ...

    Обучение на основе алгоритма Double SARSA
    ...

    self.episodes_reward = []
    # Цикл по эпизодам
    for ep in tqdm(list(range(self.num_episodes))):
        # Начальное состояние среды
        state = self.get_state(self.env.reset())
        # Флаг штатного завершения эпизода
        done = False

```

```

        # Флаг нештатного завершения эпизода
        truncated = False
        # Суммарная награда по эпизоду
        tot_rew = 0

        # По мере заполнения Q-матрицы уменьшаем вероятность случайного
        # выбора действия
        if self.eps > self.eps_threshold:
            self.eps -= self.eps_decay
        # Выбор действия

        action = self.make_action(state)

        # Проигрывание одного эпизода до финального состояния
        while not (done or truncated):

            # Выполняем шаг в среде
            next_state, rew, done, truncated, _ = self.env.step(action)

            # Выполняем следующее действие
            next_action = self.make_action(next_state)

            self.Q[state][action] = self.Q[state][action] + self.lr * \
                (rew + self.gamma * self.Q2[next_state][next_action] -
                self.Q[state][action])

            # Следующее состояние считаем текущим
            state = next_state
            action = next_action
            if np.random.rand() < 0.5:
                swap_list = self.Q
                self.Q = self.Q2
                self.Q2 = swap_list
            # Суммарная награда за эпизод
            tot_rew += rew
            if (done or truncated):
                self.episodes_reward.append(tot_rew)
def sum_rewards(self):
    # Суммарная награда
    sum_rewards = sum(self.episodes_reward)
    print('Суммарная награда Double SARSA: ', sum_rewards)

def play_agent(agent):
    """
    Проигрывание сессии для обученного агента
    """
    env2 = gym.make('CliffWalking-v0', render_mode='human') #FrozenLake-v1
    state = env2.reset()[0]
    done = False
    while not done:
        action = agent.greedy(state)

```

```

        next_state, reward, terminated, truncated, _ = env2.step(action)
        env2.render()
        state = next_state
        if terminated or truncated:
            done = True

def run_double_sarsa():
    env = gym.make('CliffWalking-v0')
    agent = DoubleSARSA_Agent(env)
    agent.learn()
    agent.sum_rewards()
    agent.print_q()
    agent.draw_episodes_reward()
    play_agent(agent)

def main():
    run_double_sarsa()

if __name__ == '__main__':
    main()

```

Результаты выполнения программы

Обучение проводилось со следующими гиперпараметрами: $\epsilon=0.015$, $lr=0.1$, $\gamma=0.97$, $\text{num_episodes}=20000$. Награды по эпизодам представлены на рис.1. Суммарная награда -333627.

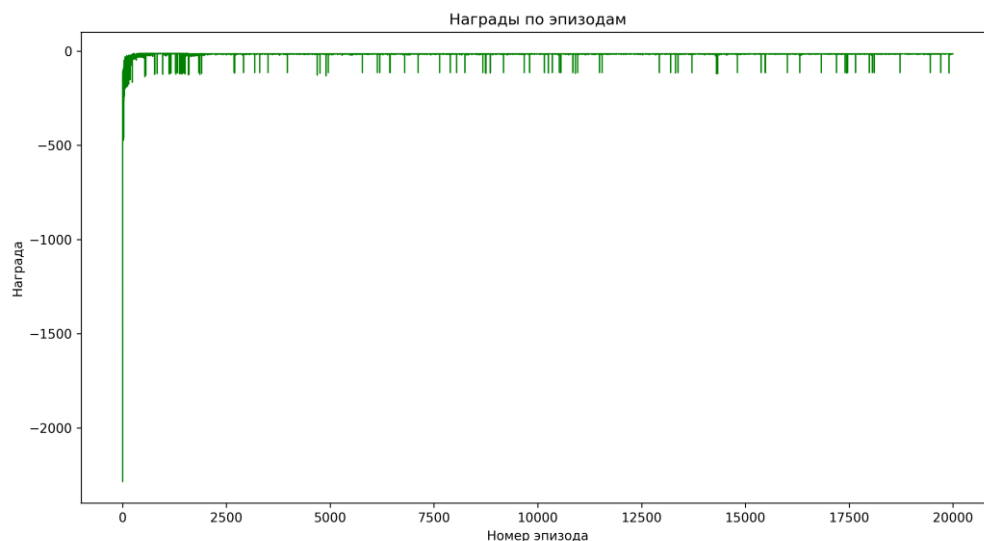


Рис.1. Награды по эпизодам

Q-матрицы для алгоритма Double SARSA представлены на рис.2 – рис.3.

Q1

[

Рис.2. Матрица Q1

состояния (клетка над конечным состоянием) соответствует исходному -1 за действие «ВНИЗ».

Путь обученного агента представлен на рис.4. При значении $\epsilon=0.4$ путь проходит по максимально «безопасной» траектории, то есть через максимально удаленные клетки от «обрыва». Однако в данном примере $\epsilon=0.015$, и стратегия дальше приближается к оптимальной при дальнейшем уменьшении этого параметра.

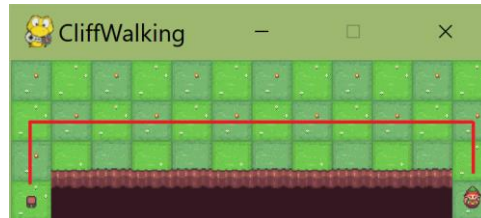


Рис.4. Путь обученного агента

Вывод: в ходе выполнения домашнего задания я ознакомилась с методом обучения с подкреплением на основе временных различий Double SARSA.

Список использованных источников

1. Ganger, M., Du-ryea, E. and Hu, W. (2016) Double Sarsa and Double Expected Sarsa with Shallow and Deep Learning. Journal of Data Analysis and Information Processing, 4, 159-176. <http://dx.doi.org/10.4236/jdaip.2016.44014>
2. Саттон Р. С., Барто Э. Дж. Обучение с подкреплением: Введение. 2-е изд. / пер. с англ. А. А. Слинкина. – М.: ДМК Пресс, 2020. – 552 с.