

Московский государственный технический университет им. Н.Э. Баумана

Факультет «Информатика, искусственный интеллект и системы управления»

Кафедра «Системы обработки информации и управления»



## Отчет по рубежному контролю №2

### по дисциплине «Методы машинного обучения»

Методы обучения с подкреплением

(тема работы)

ИСПОЛНИТЕЛЬ:

Пасатюк А.Д.

группа ИУ5-23М

ПРЕПОДАВАТЕЛЬ:

Гапанюк Ю.А.

Москва, 2023

---

## Задание

Для одного из алгоритмов временных различий, реализованных Вами в соответствующей лабораторной работе:

- SARSA
- Q-обучение
- Двойное Q-обучение

осуществите подбор гиперпараметров. Критерием оптимизации должна являться суммарная награда.

## Выполнение

Осуществим подбор гиперпараметров для алгоритма двойное Q-обучение для среды Toy Text / CliffWalking-v0.

Код 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(self):
'''
Реализация алгоритма обучения
'''
pass

# ***** Двойное Q-обучение
# *****

class DoubleQLearning_Agent(BasicAgent):
'''
Реализация алгоритма Double Q-Learning
'''
# Наименование алгоритма
ALGO_NAME = 'Двойное Q-обучение'

def __init__(self, env, eps=0.4, lr=0.01, gamma=0.98,
num_episodes=10000):
# 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 Q-Learning
    """
    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

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

            # Выбор действия
            # В SARSA следующее действие выбиралось после шага в среде
            action = self.make_action(state)

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

            if np.random.rand() < 0.5:
                # Обновление первой таблицы
                self.Q[state][action] = self.Q[state][action] + self.lr *
                \
                    (rew + self.gamma *
                     self.Q2[next_state][np.argmax(self.Q[next_state])]) -
                    self.Q[state][action])
            else:

```

```

        # Обновление второй таблицы
        self.Q2[state][action] = self.Q2[state][action] + self.lr
* \
        (rew + self.gamma *
self.Q[next_state][np.argmax(self.Q2[next_state])] -
        self.Q2[state][action])

        # Следующее состояние считаем текущим
        state = next_state
        # Суммарная награда за эпизод
        tot_rew += rew
        if (done or truncated):
            self.episodes_reward.append(tot_rew)

def play_agent(agent):
    """
    Проигрывание сессии для обученного агента
    """
    env2 = gym.make('CliffWalking-v0', render_mode='human')
    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_sarsa():
    env = gym.make('CliffWalking-v0')
    agent = SARSA_Agent(env)
    agent.learn()
    agent.print_q()
    agent.draw_episodes_reward()
    play_agent(agent)

def run_q_learning():
    env = gym.make('CliffWalking-v0')
    agent = QLearning_Agent(env)
    agent.learn()
    agent.print_q()
    agent.draw_episodes_reward()
    play_agent(agent)

def run_double_q_learning():
    env = gym.make('CliffWalking-v0')
    agent = DoubleQLearning_Agent(env)
    agent.learn()
    agent.print_q()
    agent.draw_episodes_reward()
    play_agent(agent)

def main():
    #run_sarsa()
    #run_q_learning()
    run_double_q_learning()

```

```
if __name__ == '__main__':  
    main()
```

Начальные значения параметров: `eps=0.4`, `lr=0.01`, `gamma=0.98`,  
`num_episodes=20000`

Результат работы программы для алгоритма двойное Q-обучение:

Суммарная награда: -1 364 805

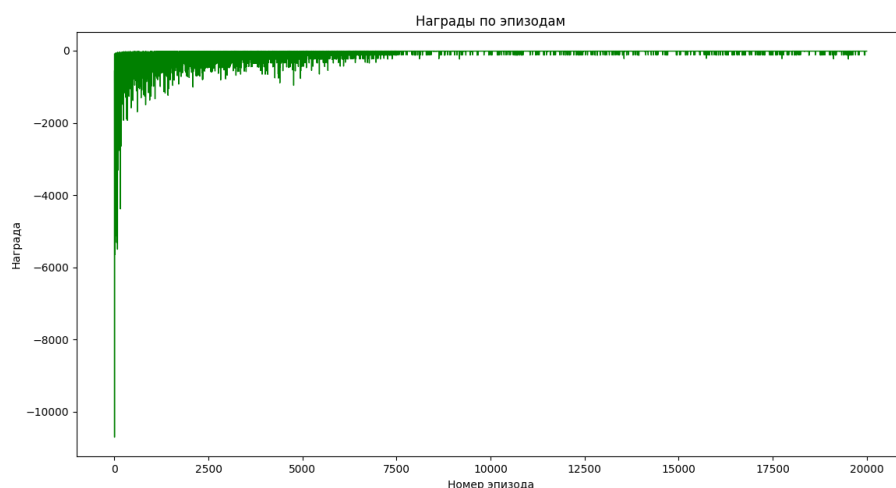
Вывод Q-матриц для алгоритма Двойное Q-обучение

Q1

[illegible]

02

```
[[ -10.08725201  -10.12520674  -10.1195806   -10.11902158]
```

[illegible]

Изменим следующие параметры: lr=0.5, num\_episodes=20000

Суммарная награда: -1 439 135

Вывод Q-матриц для алгоритма Двойное Q-обучение

Q1

[	-11.9657858	-11.87042023	-11.86732539	-11.83872938]
[	-11.41332339	-11.39771521	-11.4175467	-11.74805093]
[	-10.68288564	-10.71451528	-10.67901651	-10.91284192]
[	-9.89879522	-9.86997088	-9.93874431	-10.33062944]
[	-9.51627825	-9.23844394	-9.13758822	-9.6380229 ]
[	-8.48133609	-8.5216577	-8.28248331	-8.823221 ]
[	-7.61497493	-7.87381806	-8.1313079	-7.66800254]
[	-7.23893188	-6.72829976	-6.41499778	-7.0540342 ]
[	-5.40910339	-6.24544729	-6.52080154	-5.86225381]
[	-5.34673589	-4.73160967	-4.20692848	-4.83008464]
[	-3.90981407	-3.925385	-4.52694793	-4.34713804]
[	-3.44020774	-3.27828083	-2.93761498	-3.54764939]
[	-12.46313118	-11.61054848	-11.54888054	-12.19742999]
[	-11.91560349	-10.75625267	-10.76416381	-12.22470474]
[	-11.23176454	-9.97286426	-9.96343246	-11.53577846]
[	-10.55222519	-9.17263827	-9.14635966	-10.6015 ]
[	-9.8356837	-8.31490101	-8.31261189	-9.80883516]
[	-9.44160891	-8.18802267	-7.46184887	-8.93985243]
[	-8.25889616	-6.52998622	-6.59372333	-8.07883231]
[	-7.77040861	-6.76957049	-5.70788099	-7.66870773]
[	-6.23967876	-4.45823857	-4.80395985	-6.14192221]
[	-5.92363029	-4.77522461	-3.88159231	-6.30472527]
[	-4.25804498	-2.93969628	-2.93976962	-3.72663587]
[	-3.55458229	-2.72681421	-1.98	-3.9549554 ]
[	-12.31790293	-10.76416381	-12.31790293	-11.54888054]
[	-11.57873915	-9.96343246	-111.31790293	-11.54888054]
[	-10.76416381	-9.14635966	-111.31790293	-10.76416381]
[	-9.96343246	-8.31261189	-111.31790293	-9.96343246]
[	-9.14635966	-7.46184887	-111.31790293	-9.14635966]
[	-8.31261189	-6.59372334	-111.3179029	-8.31261189]
[	-8.31503634	-5.70788096	-111.31790282	-7.46184887]
[	-6.59372283	-4.80396016	-111.31790276	-6.59372322]
[	-6.80705805	-3.881592	-111.31789958	-5.70788098]
[	-4.80395568	-2.9404	-111.31789651	-4.80395934]
[	-3.88234228	-1.98	-111.31788262	-3.88159162]
[	-2.94039874	-1.97999827	-1.	-2.94039945]
[	-11.54888054	-111.31790293	-12.31790293	-12.31790293]
[	0.	0.	0.	0. ]
[	0.	0.	0.	0. ]
[	0.	0.	0.	0. ]
[	0.	0.	0.	0. ]
[	0.	0.	0.	0. ]
[	0.	0.	0.	0. ]
[	0.	0.	0.	0. ]
[	0.	0.	0.	0. ]
[	0.	0.	0.	0. ]
[	0.	0.	0.	0. ]
[	0.	0.	0.	0. ]

Q2

[	-11.85627492	-11.94110353	-11.9363446	-11.99527241]
[	-11.36027433	-11.37766876	-11.36031981	-11.43012568]
[	-10.86007128	-10.67902143	-10.71419049	-10.98906902]
[	-10.1923441	-9.98868682	-9.91805686	-10.12927801]
[	-9.27187704	-9.12713996	-9.14280062	-9.41085218]
[	-8.65797006	-8.60423479	-8.82751155	-8.91765225]
[	-7.97369461	-7.55849936	-7.27363464	-8.16736263]
[	-6.60225054	-7.0606398	-7.38036484	-7.18632466]





Изменим следующие параметры: `lr=0.8`, `num_episodes=10000`

Суммарная награда: -859 934

Вывод Q-матриц для алгоритма Двойное Q-обучение

Q1

```
[[ -13.07239812 -12.31790978 -12.31790293 -13.07240228]
 [ -12.31790293 -11.54888056 -11.54888054 -13.07154487]
 [ -11.54888059 -10.76416381 -10.76416381 -12.31790294]
 [ -10.76416381 -9.96346368 -9.96343246 -11.54888054]
 [ -9.96343246 -9.14635966 -9.14635966 -10.76416381]
 [ -9.14635966 -8.31261189 -8.31261189 -9.96343472]
 [ -8.31261189 -7.46184887 -7.46184887 -9.14635966]
 [ -7.46184887 -6.59372334 -6.59372334 -8.31261189]
 [ -6.59372334 -5.70788099 -5.70788096 -7.46184887]
 [ -5.70788096 -4.80396016 -4.80396016 -6.59372334]
 [ -4.80396016 -3.881592 -3.881592 -5.70788096]
 [ -3.881592 -3.881592 -2.9404 -4.80396016]
 [ -13.07154487 -11.54888054 -11.54888054 -12.31790293]
 [ -12.31790293 -10.76416381 -10.76416381 -12.31790293]
 [ -11.54888054 -9.96343246 -9.96343246 -11.54888054]
 [ -10.76416381 -9.14635966 -9.14635966 -10.76416381]
 [ -9.96343246 -8.31261189 -8.31261189 -9.96343246]
 [ -9.14635966 -7.46184887 -7.46184887 -9.14635966]
 [ -8.31261189 -6.59372334 -6.59372334 -8.31261189]
 [ -7.46184887 -5.70788096 -5.70788096 -7.46184887]
 [ -6.59372334 -4.80396016 -4.80396016 -6.59372334]
 [ -5.70788096 -3.881592 -3.881592 -5.70788096]
 [ -4.80396016 -2.9404 -2.9404 -4.80396016]
 [ -3.881592 -2.9404 -1.98 -3.881592 ]
 [ -12.31790293 -10.76416381 -12.31790293 -11.54888054]
 [ -11.54888054 -9.96343246 -11.31790293 -11.54888054]
 [ -10.76416381 -9.14635966 -11.31790293 -10.76416381]
 [ -9.96343246 -8.31261189 -11.31790293 -9.96343246]
 [ -9.14635966 -7.46184887 -11.31790293 -9.14635966]
 [ -8.31261189 -6.59372334 -11.31790293 -8.31261189]
 [ -7.46184887 -5.70788096 -11.31790293 -7.46184887]
 [ -6.59372334 -4.80396016 -11.31790293 -6.59372334]
 [ -5.70788096 -3.881592 -11.31790293 -5.70788096]
 [ -4.80396016 -2.9404 -11.31790293 -4.80396016]
 [ -3.881592 -1.98 -11.31790293 -3.881592 ]
 [ -2.9404 -1.98 -1. -2.9404 ]
 [ -11.54888054 -11.31790293 -12.31790293 -12.31790293]
 [ 0. 0. 0. 0. ]
 [ 0. 0. 0. 0. ]
 [ 0. 0. 0. 0. ]
 [ 0. 0. 0. 0. ]
 [ 0. 0. 0. 0. ]
 [ 0. 0. 0. 0. ]
 [ 0. 0. 0. 0. ]
 [ 0. 0. 0. 0. ]
 [ 0. 0. 0. 0. ]
 [ 0. 0. 0. 0. ]
 [ 0. 0. 0. 0. ]
 [ 0. 0. 0. 0. ]]
```

Q2

```
[ [-13.07154492 -12.31790433 -12.31790293 -13.07154492]
[ -12.31790294 -11.54888054 -11.54888054 -13.07154488]
[ -11.54888054 -10.76416381 -10.76416381 -12.31790293]
[ -10.76416381 -9.96343873 -9.96343246 -11.5488808 ]
[ -9.96343246 -9.14635966 -9.14635966 -10.76416381]
[ -9.14635966 -8.3126119 -8.31261189 -9.96343246]
[ -8.3126119 -7.46184887 -7.46184888 -9.14636125]
[ -7.46184887 -6.59372334 -6.59372334 -8.31261189]
[ -6.59372334 -5.70788125 -5.70788096 -7.46184887]
[ -5.70788099 -4.80396016 -4.80396016 -6.59372337]
[ -4.80396016 -3.881592 -3.881592 -5.70788096]
[ -3.881592 -3.881592 -2.9404 -4.80396016]
[ -13.07154487 -11.54888054 -11.54888054 -12.31790293]
[ -12.31790293 -10.76416381 -10.76416381 -12.31790293]
[ -11.54888054 -9.96343246 -9.96343246 -11.54888054]
[ -10.76416381 -9.14635966 -9.14635966 -10.76416381]
[ -9.96343246 -8.31261189 -8.31261189 -9.96343246]
[ -9.14635966 -7.46184887 -7.46184887 -9.14635966]
[ -8.31261189 -6.59372334 -6.59372334 -8.31261189]
[ -7.46184887 -5.70788096 -5.70788096 -7.46184887]
[ -6.59372334 -4.80396016 -4.80396016 -6.59372334]
[ -5.70788096 -3.881592 -3.881592 -5.70788096]
[ -4.80396016 -2.9404 -2.9404 -4.80396016]
[ -3.881592 -2.9404 -1.98 -3.881592 ]
[ -12.31790293 -10.76416381 -12.31790293 -11.54888054]
[ -11.54888054 -9.96343246 -11.31790293 -11.54888054]
[ -10.76416381 -9.14635966 -11.31790293 -10.76416381]
[ -9.96343246 -8.31261189 -11.31790293 -9.96343246]
[ -9.14635966 -7.46184887 -11.31790293 -9.14635966]
[ -8.31261189 -6.59372334 -11.31790293 -8.31261189]
[ -7.46184887 -5.70788096 -11.31790293 -7.46184887]
[ -6.59372334 -4.80396016 -11.31790293 -6.59372334]
[ -5.70788096 -3.881592 -11.31790293 -5.70788096]
[ -4.80396016 -2.9404 -11.31790293 -4.80396016]
[ -3.881592 -1.98 -11.31790293 -3.881592 ]
[ -2.9404 -1.98 -1. -2.9404 ]
[ -11.54888054 -11.31790293 -12.31790293 -12.31790293]]
```

