



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

ФАКУЛЬТЕТ Информатика и системы управления и искусственный интеллект

КАФЕДРА Системы обработки информации и управления

Лабораторная работа №5
По курсу
«Методы машинного обучения в АСОИУ»
«Обучение на основе временных различий»

Выполнил:

ИУ5-22М Киричков Е. Е.

22.05.2024

Проверил:

Балашов А.М.

Москва, 2024

Задание:

На основе рассмотренного на лекции примера реализуйте следующие алгоритмы:

- SARSA
- Q-обучение
- Двойное Q-обучение

для любой среды обучения с подкреплением (кроме рассмотренной на лекции среды Toy Text / Frozen Lake) из библиотеки Gym (или аналогичной библиотеки).

```
Ввод [1]: import numpy as np
import matplotlib.pyplot as plt
import gym
from tqdm import tqdm
from IPython.display import Image
```

Ввод [2]: # ***** 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 но
            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
```

Ввод [3]: # ***** SARSA *****

```
class SARSA_Agent(BasicAgent):
    """
    Реализация алгоритма SARSA
    """
    # Наименование алгоритма
    ALGO_NAME = 'SARSA'

    def __init__(self, env, eps=0.4, lr=0.1, gamma=0.98, num_episodes=20000):
        # Вызов конструктора верхнего уровня
        super().__init__(env, eps)
        # 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 learn(self):
        """
        Обучение на основе алгоритма 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
            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)

                # Правило обновления Q для SARSA
                self.Q[state][action] = self.Q[state][action] + self.lr * \
                    (rew + self.gamma * self.Q[next_state][next_action] - s

                # Следующее состояние считаем текущим
                state = next_state
                action = next_action
```

```
# Суммарная награда за эпизод
tot_rew += rew
if (done or truncated):
    self.episodes_reward.append(tot_rew)
```

Ввод [4]: # ***** Q-обучение *****

```
class QLearning_Agent(BasicAgent):
    """
    Реализация алгоритма Q-Learning
    """
    # Наименование алгоритма
    ALGO_NAME = 'Q-обучение'

    def __init__(self, env, eps=0.4, lr=0.1, gamma=0.98, num_episodes=20000):
        # Вызов конструктора верхнего уровня
        super().__init__(env, eps)
        # 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 learn(self):
        """
        Обучение на основе алгоритма 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
            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)

                # Правило обновления Q для SARSA (для сравнения)
                # self.Q[state][action] = self.Q[state][action] + self.lr * \
                #     (rew + self.gamma * self.Q[next_state][next_action] -

                # Правило обновления для Q-обучения
                self.Q[state][action] = self.Q[state][action] + self.lr * \
                    (rew + self.gamma * np.max(self.Q[next_state]) - self.Q

                # Следующее состояние считаем текущим
```

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


Ввод [5]: # ***** 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=0.1, 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 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
            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(s
    else:
        # Обновление второй таблицы
        self.Q2[state][action] = self.Q2[state][action] + self.
            (rew + self.gamma * self.Q[next_state][np.argmax(se

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

```

Ввод [6]:

```

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

```

```
Ввод [7]: def run_sarsa():
            env = gym.make('CliffWalking-v0')
            agent = SARSA_Agent(env)
            agent.learn()
            agent.print_q()
            agent.draw_episodes_reward()
            play_agent(agent)

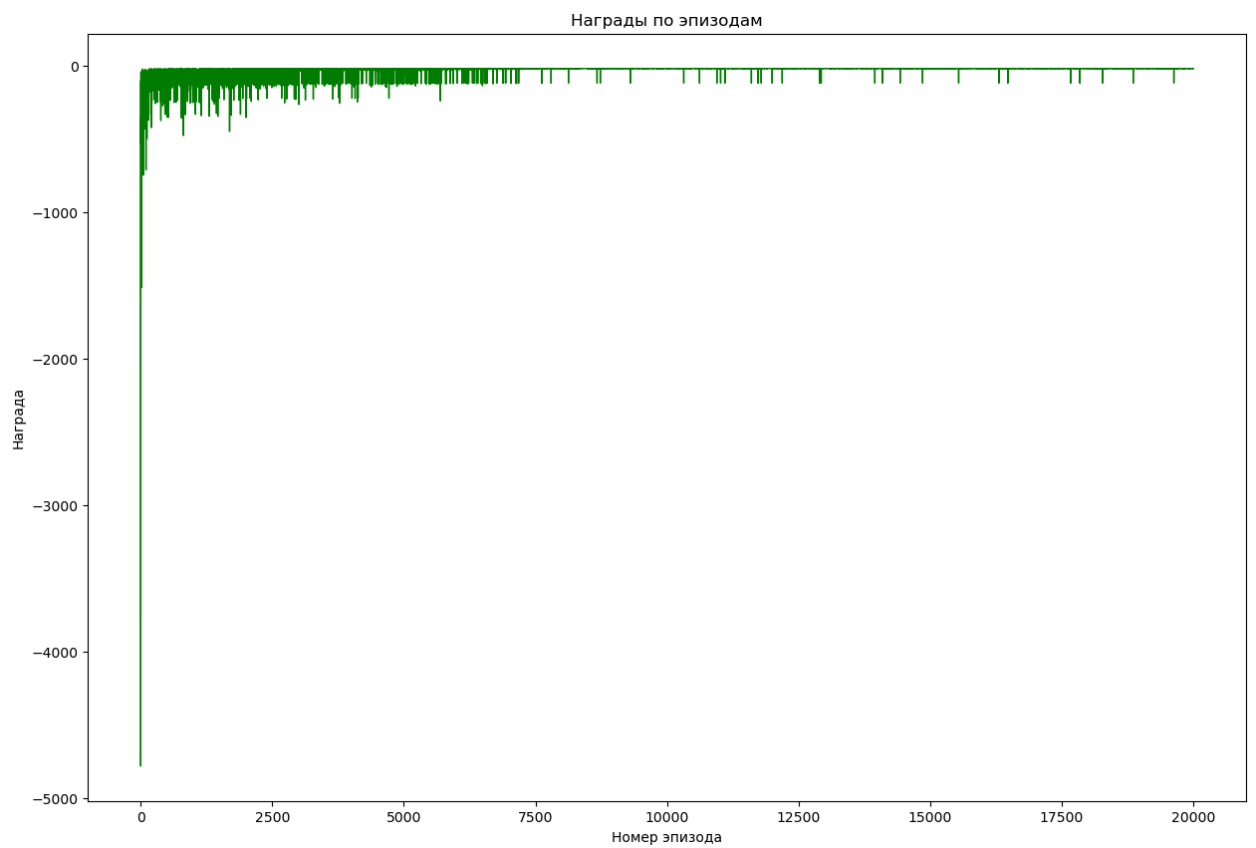
def run_q_learning():
    env = gym.make('CliffWalking-v0')
    agent = QLearning_Agent(env)
    agent.learn()
    agent.print_q()
    agent.draw_episodes_reward()
    play_agent(agent)

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)
```

```
Ввод [8]: def main():
            run_sarsa()
            # run_q_learning()
            # run_double_q_learning()
```

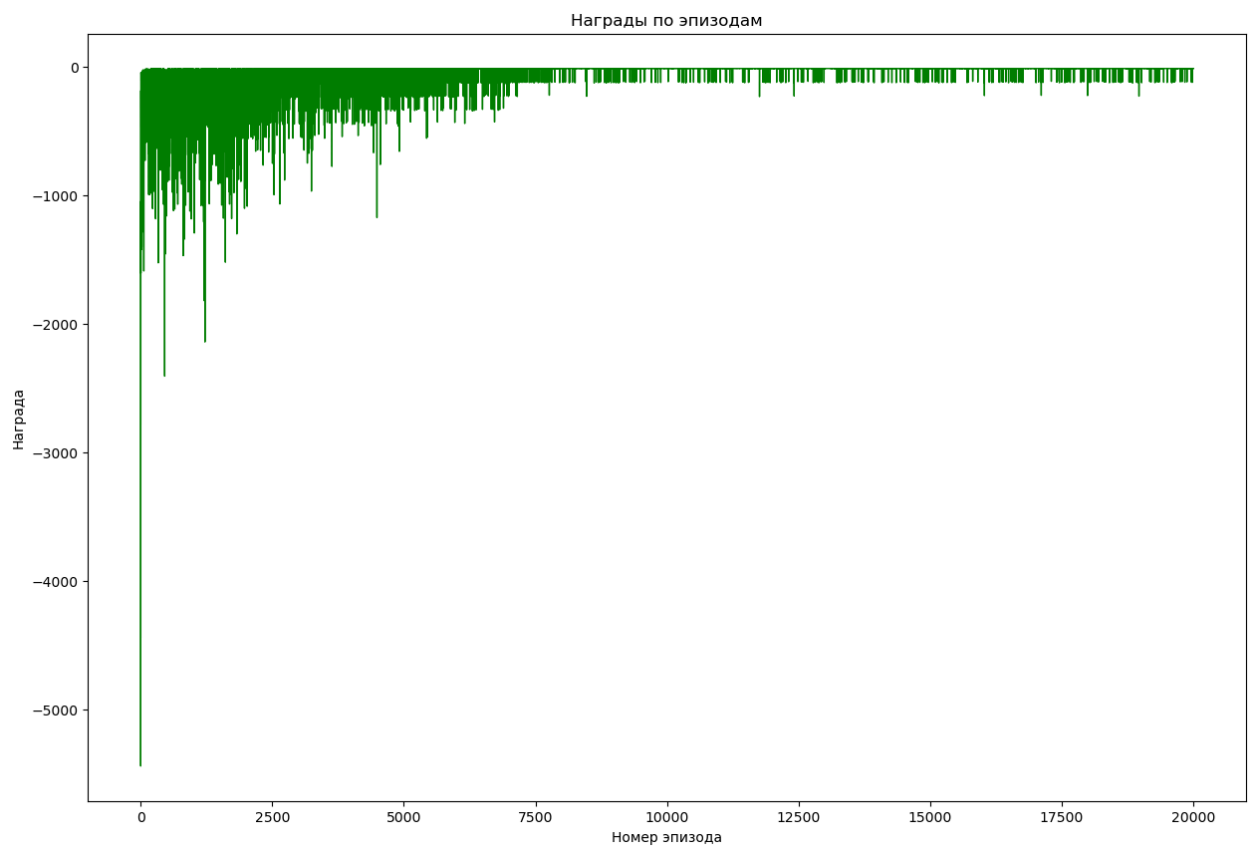
```
Ввод [9]: if __name__ == '__main__':  
           main()
```

[illegible][illegible]



```
Ввод [10]: image_path_1 = "data/sarsa_result.png"  
display(Image(filename=image_path_1, width=600, height=600))
```





```
Ввод [12]: image_path_2 = "data/q_learning_result.png"  
display(Image(filename=image_path_2, width=600, height=600))
```



```
Ввод [13]: run_double_q_learning()
```

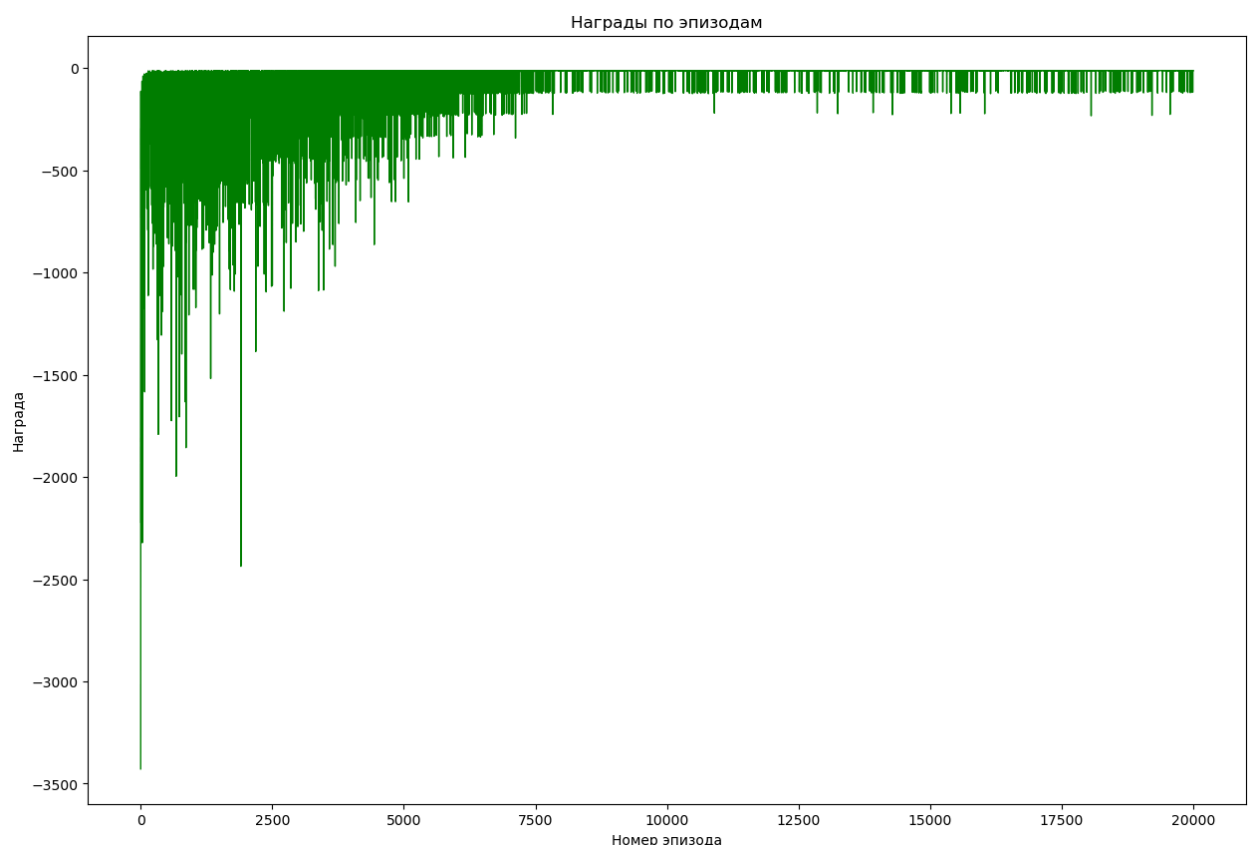
```
100%|██████████████████████████████████████████████████████████████████████████  
██████████████████████████████████████████████████████████████████████████████| 20000/20000 [00:02<0  
0:00, 6837.80it/s]
```


Q1

Q2

```
[ [ -12.70448892 -12.73304683 -12.32201258 -13.00957706]
  [ -12.46269169 -11.5613982 -11.89280967 -12.55332763]
  [ -11.54122143 -10.96027511 -10.76417181 -12.2294357 ]
  [ -10.83487034 -9.96344275 -10.15643485 -11.42368401]
  [ -9.99140083 -9.43053819 -9.14635966 -10.66587436]
  [ -9.29760507 -8.39053038 -8.31261222 -9.85891361]
  [ -8.12678046 -7.71423862 -7.46533718 -8.8435014 ]
  [ -7.55817293 -6.53971499 -6.59372333 -8.22062163]
  [ -6.13661645 -5.8543554 -5.79124268 -7.02860736]
  [ -5.268676 -4.73183315 -4.8036398 -6.06501667]
  [ -4.33076759 -3.93296974 -3.88100595 -5.0057501 ]
  [ -3.49841432 -3.05033065 -2.94101726 -4.02710949]
```

[-13.08759308	-11.55833624	-11.54888054	-12.3297932]
[-12.37343169	-10.76416381	-10.78169256	-12.3241107]
[-11.55054966	-9.96343246	-9.96362363	-11.54925377]
[-10.76528321	-9.14635966	-9.14643843	-10.76423969]
[-9.96350015	-8.31261189	-8.31262288	-9.96344855]
[-9.14643258	-7.46184887	-7.46186527	-9.14636147]
[-8.29130484	-6.59372334	-6.59372569	-8.31260495]
[-7.68058395	-5.74006026	-5.70788096	-7.46185269]
[-6.47095046	-4.80396016	-4.78870442	-6.49809894]
[-5.83999245	-3.881592	-3.88169182	-5.77271787]
[-4.80036846	-2.94055416	-2.9404	-4.80071899]
[-3.93683849	-2.92169017	-1.98	-3.79101561]
[-12.31790293	-10.76416381	-12.31790293	-11.54888054]
[-11.54888054	-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	-111.31790293	-8.31261189]
[-7.46184887	-5.70788096	-111.31790293	-7.46184887]
[-6.59372334	-4.80396016	-111.31790293	-6.59372334]
[-5.72545506	-3.881592	-111.31790293	-5.70788096]
[-4.80396016	-2.9404	-111.31790293	-4.80396016]
[-3.881592	-1.98	-111.31790293	-3.881592]
[-2.9404	-1.98	-1.	-2.9404]
[-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.]
[0.	0.	0.	0.]
[0.	0.	0.	0.]



```
Ввод [14]: image_path_3 = "data/double_q_learning_result.png"
display(Image(filename=image_path_3, width=600, height=600))
```



Итог

SARSA учится, основываясь на том, что агент действительно делает, включая случайные действия. Это делает путь агента более осторожным и безопасным, потому что алгоритм учитывает риск этих случайных действий. Q-learning учится на основе лучшего возможного действия в каждой ситуации, не учитывая случайные действия. Это делает путь агента более прямым и рискованным, потому что он всегда стремится к максимальной награде, даже если это может быть опасно. Double Q-learning направлен на устранение смещения переоценки в Q-learning, поддерживая два отдельных Q-значения и обновляя их поочередно. Это обычно приводит к более сбалансированному и стабильному процессу обучения. Путь, выбранный агентом с использованием Double Q-learning, обычно является компромиссом между осторожным подходом SARSA и агрессивным подходом Q-learning, что приводит к более надежной политике.

Ввод []: