



Федеральное государственное бюджетное образовательное учреждение
высшего профессионального образования
«Московский государственный технический университет имени Н.Э.
Баумана»
(МГТУ им. Н.Э. Баумана)

**Лабораторная работа №5
по курсу «Методы машинного обучения»**

Выполнил
студент группы ИУ5-22М
XXXX

Москва, 2023

1. Задание

На основе рассмотренного на лекции примера реализуйте следующие алгоритмы:

- SARSA
- Q-обучение
- Двойное Q-обучение

2. Текст программы

```
1  #!/usr/bin/env python
2
3  import numpy as np
4  import matplotlib.pyplot as plt
5  import gymnasium as gym
6  from tqdm import tqdm
7
8  # ***** БАЗОВЫЙ АГЕНТ *****
9
10 class BasicAgent:
11     '''
12     Базовый агент, от которого наследуются стратегии обучения
13     '''
14
15     # Наименование алгоритма
16     ALGO_NAME = '___'
17
18     def __init__(self, env, eps=0.1):
19         # Среда
20         self.env = env
21         # Размерности Q-матрицы
22         self.nA = env.action_space.n
23         self.nS = env.observation_space.n
24         #и сама матрица
25         self.Q = np.zeros((self.nS, self.nA))
26         # Значения коэффициентов
27         # Порог выбора случайного действия
28         self.eps=eps
29         # Награды по эпизодам
30         self.episodes_reward = []
31
32     def print_q(self):
33         print('Вывод Q-матрицы для алгоритма ', self.ALGO_NAME)
34         print(self.Q)
35
36     def get_state(self, state):
37         '''
38         Возвращает правильное начальное состояние
39         '''
40         if type(state) is tuple:
41             # Если состояние вернулось с виде кортежа, то вернуть только номер состояния
42             return state[0]
43         else:
44             return state
45
46     def greedy(self, state):
```

```

47     '''
48     <<Жадное>> текущее действие
49     Возвращает действие, соответствующее максимальному Q-значению
50     для состояния state
51     '''
52     return np.argmax(self.Q[state])
53
54 def make_action(self, state):
55     '''
56     Выбор действия агентом
57     '''
58     if np.random.uniform(0,1) < self.eps:
59
60         # Если вероятность меньше eps
61         # то выбирается случайное действие
62         return self.env.action_space.sample()
63     else:
64         # иначе действие, соответствующее максимальному Q-значению
65         return self.greedy(state)
66
67 def draw_episodes_reward(self):
68     # Построение графика наград по эпизодам
69     fig, ax = plt.subplots(figsize = (15,10))
70     y = self.episodes_reward
71     x = list(range(1, len(y)+1))
72     plt.plot(x, y, '-', linewidth=1, color='green')
73     plt.title('Награды по эпизодам')
74     plt.xlabel('Номер эпизода')
75     plt.ylabel('Награда')
76     plt.show()
77
78 def learn():
79     '''
80     Реализация алгоритма обучения
81     '''
82     pass
83
84 # ***** SARSA *****
85
86 class SARSA_Agent(BasicAgent):
87     '''
88     Реализация алгоритма SARSA
89     '''
90     # Наименование алгоритма
91     ALGO_NAME = 'SARSA'
92
93     def __init__(self, env, eps=0.4, lr=0.1, gamma=0.98, num_episodes=20000):
94         # Вызов конструктора верхнего уровня
95         super().__init__(env, eps)
96         # Learning rate
97         self.lr=lr
98         # Коэффициент дисконтирования
99         self.gamma = gamma
100        # Количество эпизодов
101        self.num_episodes=num_episodes
102        # Постепенное уменьшение eps
103        self.eps_decay=0.00005
104        self.eps_threshold=0.01
105
106    def learn(self):
107        '''

```

```

108     Обучение на основе алгоритма SARSA
109     '''
110     self.episodes_reward = []
111     # Цикл по эпизодам
112     for ep in tqdm(list(range(self.num_episodes))):
113         # Начальное состояние среды
114         state = self.get_state(self.env.reset())
115         # Флаг штатного завершения эпизода
116         done = False
117         # Флаг нештатного завершения эпизода
118         truncated = False
119         # Суммарная награда по эпизоду
120         tot_rew = 0
121
122         # По мере заполнения Q-матрицы уменьшаем вероятность случайного выбора действия
123         if self.eps > self.eps_threshold:
124             self.eps -= self.eps_decay
125
126         # Выбор действия
127         action = self.make_action(state)
128
129         # Проигрывание одного эпизода до финального состояния
130         while not (done or truncated):
131             # Выполняем шаг в среде
132             next_state, rew, done, truncated, _ = self.env.step(action)
133
134             # Выполняем следующее действие
135             next_action = self.make_action(next_state)
136
137             # Правило обновления Q для SARSA
138             self.Q[state][action] = self.Q[state][action] + self.lr * \
139                 (rew + self.gamma * self.Q[next_state][next_action] -
140                  self.Q[state][action])
141
142             # Следующее состояние считаем текущим
143             state = next_state
144             action = next_action
145             # Суммарная награда за эпизод
146             tot_rew += rew
147             if (done or truncated):
148                 self.episodes_reward.append(tot_rew)
149 # ***** Q-обучение *****
150
151 class QLearning_Agent(BasicAgent):
152     '''
153     Реализация алгоритма Q-Learning
154     '''
155     # Наименование алгоритма
156     ALGO_NAME = 'Q-обучение'
157
158     def __init__(self, env, eps=0.4, lr=0.1, gamma=0.98, num_episodes=20000):
159         # Вызов конструктора верхнего уровня
160         super().__init__(env, eps)
161         # Learning rate
162         self.lr=lr
163         # Коэффициент дисконтирования
164         self.gamma = gamma
165         # Количество эпизодов
166         self.num_episodes=num_episodes
167         # Постепенное уменьшение eps

```

```

168         self.eps_decay=0.00005
169         self.eps_threshold=0.01
170
171     def learn(self):
172         '''
173         Обучение на основе алгоритма Q-Learning
174         '''
175         self.episodes_reward = []
176         # Цикл по эпизодам
177         for ep in tqdm(list(range(self.num_episodes))):
178             # Начальное состояние среды
179             state = self.get_state(self.env.reset())
180             # Флаг штатного завершения эпизода
181             done = False
182             # Флаг нештатного завершения эпизода
183             truncated = False
184             # Суммарная награда по эпизоду
185             tot_rew = 0
186
187             # По мере заполнения Q-матрицы уменьшаем вероятность случайного выбора действия
188             if self.eps > self.eps_threshold:
189                 self.eps -= self.eps_decay
190
191             # Проигрывание одного эпизода до финального состояния
192             while not (done or truncated):
193                 # Выбор действия
194                 # В SARSA следующее действие выбиралось после шага в среде
195                 action = self.make_action(state)
196
197                 # Выполняем шаг в среде
198                 next_state, rew, done, truncated, _ = self.env.step(action)
199
200                 # Правило обновления Q для SARSA (для сравнения)
201                 # self.Q[state][action] = self.Q[state][action] + self.lr * \
202                 #     (rew + self.gamma * self.Q[next_state][next_action] -
203                 #      self.Q[state][action])
204
205                 # Правило обновления для Q-обучения
206                 self.Q[state][action] = self.Q[state][action] + self.lr * \
207                     (rew + self.gamma * np.max(self.Q[next_state]) - self.Q[state][action])
208
209                 # Следующее состояние считаем текущим
210                 state = next_state
211                 # Суммарная награда за эпизод
212                 tot_rew += rew
213                 if (done or truncated):
214                     self.episodes_reward.append(tot_rew)
215
216     # ***** Двойное Q-обучение *****
217
218 class DoubleQLearning_Agent(BasicAgent):
219     '''
220     Реализация алгоритма Double Q-Learning
221     '''
222     # Наименование алгоритма
223     ALGO_NAME = 'Двойное Q-обучение'
224
225     def __init__(self, env, eps=0.4, lr=0.1, gamma=0.98, num_episodes=20000):
226         # Вызов конструктора верхнего уровня
227         super().__init__(env, eps)

```

```

227     # Вторая матрица
228     self.Q2 = np.zeros((self.nS, self.nA))
229     # Learning rate
230     self.lr=lr
231     # Коэффициент дисконтирования
232     self.gamma = gamma
233     # Количество эпизодов
234     self.num_episodes=num_episodes
235     # Постепенное уменьшение eps
236     self.eps_decay=0.00005
237     self.eps_threshold=0.01
238
239
240     def greedy(self, state):
241         '''
242         <<Жадное>> текущее действие
243         Возвращает действие, соответствующее максимальному Q-значению
244         для состояния state
245         '''
246         temp_q = self.Q[state] + self.Q2[state]
247         return np.argmax(temp_q)
248
249     def print_q(self):
250         print(f"Вывод Q-матриц для алгоритма {self.ALGO_NAME}")
251         print('Q1')
252         print(self.Q)
253         print('Q2')
254         print(self.Q2)
255
256     def learn(self):
257         '''
258         Обучение на основе алгоритма Double Q-Learning
259         '''
260         self.episodes_reward = []
261         # Цикл по эпизодам
262         for ep in tqdm(list(range(self.num_episodes))):
263             # Начальное состояние среды
264             state = self.get_state(self.env.reset())
265             # Флаг штатного завершения эпизода
266             done = False
267             # Флаг нештатного завершения эпизода
268             truncated = False
269             # Суммарная награда по эпизоду
270             tot_rew = 0
271
272             # По мере заполнения Q-матрицы уменьшаем вероятность случайного выбора действия
273             if self.eps > self.eps_threshold:
274                 self.eps -= self.eps_decay
275
276             # Проигрывание одного эпизода до финального состояния
277             while not (done or truncated):
278                 # Выбор действия
279                 # В SARSA следующее действие выбиралось после шага в среде
280                 action = self.make_action(state)
281
282                 # Выполняем шаг в среде
283                 next_state, rew, done, truncated, _ = self.env.step(action)
284
285                 if np.random.rand() < 0.5:
286                     # Обновление первой таблицы
287                     self.Q[state][action] = self.Q[state][action] + self.lr * \

```

```

288         (rew + self.gamma * self.Q2[next_state][np.argmax(self.Q[next_state])])
289         - self.Q[state][action])
290     else:
291         # Обновление второй таблицы
292         self.Q2[state][action] = self.Q2[state][action] + self.lr * \
293             (rew + self.gamma * self.Q2[next_state][np.argmax(self.Q2[next_state])])
294             - self.Q2[state][action])
295
296     # Следующее состояние считаем текущим
297     state = next_state
298     # Суммарная награда за эпизод
299     tot_rew += rew
300     if (done or truncated):
301         self.episodes_reward.append(tot_rew)
302
303 def play_agent(agent):
304     '''
305     Проиhrывание сессии для обученного агента
306     '''
307     env2 = gym.make('Taxi-v3', render_mode='human')
308     state = env2.reset()[0]
309     done = False
310     while not done:
311         action = agent.greedy(state)
312         next_state, reward, terminated, truncated, _ = env2.step(action)
313         env2.render()
314         state = next_state
315         if terminated or truncated:
316             done = True
317
318 def run_sarsa():
319     env = gym.make('Taxi-v3')
320     agent = SARSA_Agent(env)
321     agent.learn()
322     agent.print_q()
323     agent.draw_episodes_reward()
324     play_agent(agent)
325
326 def run_q_learning():
327     env = gym.make('Taxi-v3')
328     agent = QLearning_Agent(env)
329     agent.learn()
330     agent.print_q()
331     agent.draw_episodes_reward()
332     play_agent(agent)
333
334 def run_double_q_learning():
335     env = gym.make('Taxi-v3')
336     agent = DoubleQLearning_Agent(env)
337     agent.learn()
338     agent.print_q()
339     agent.draw_episodes_reward()
340     play_agent(agent)
341
342 def main():
343     # run_sarsa()
344     run_q_learning()
345     # run_double_q_learning()
346
347 if __name__ == '__main__':
348     main()

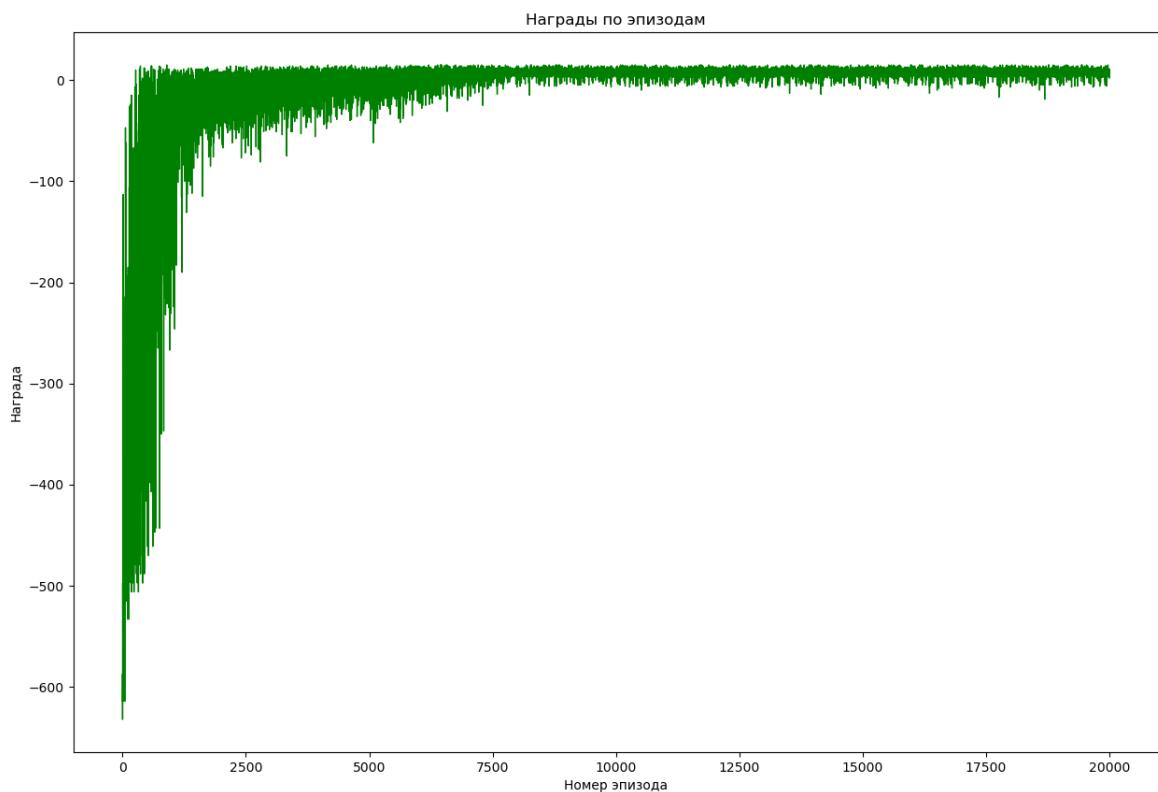
```

3. Экранные формы с примерами выполнения программы



Q-обучение: $\text{eps}=0.4$, $\text{lr}=0.1$, $\text{gamma}=0.98$, $\text{num_episodes}=20000$

```
1 Вывод Q-матрицы для алгоритма Q-обучение
2 [[ 0.          0.          0.          0.          0.          0.        ]
3  [ 5.09940217  5.11613495  4.9333549   6.24799219  8.36234335 -3.62258081]
4  [10.04299356 11.11574742  9.6393241   11.43318009 13.27445578  2.65833957]
5  ...
6  [-1.36364673 11.89866367  0.31299649  0.12111493 -2.18307592 -4.64070247]
7  [-2.6539698  8.34882836 -2.02398777  0.61115667 -8.4092672  -7.78482003]
8  [11.04932363  4.15528476  7.05358754 18.59990519  2.11145862  4.24396103]]
```

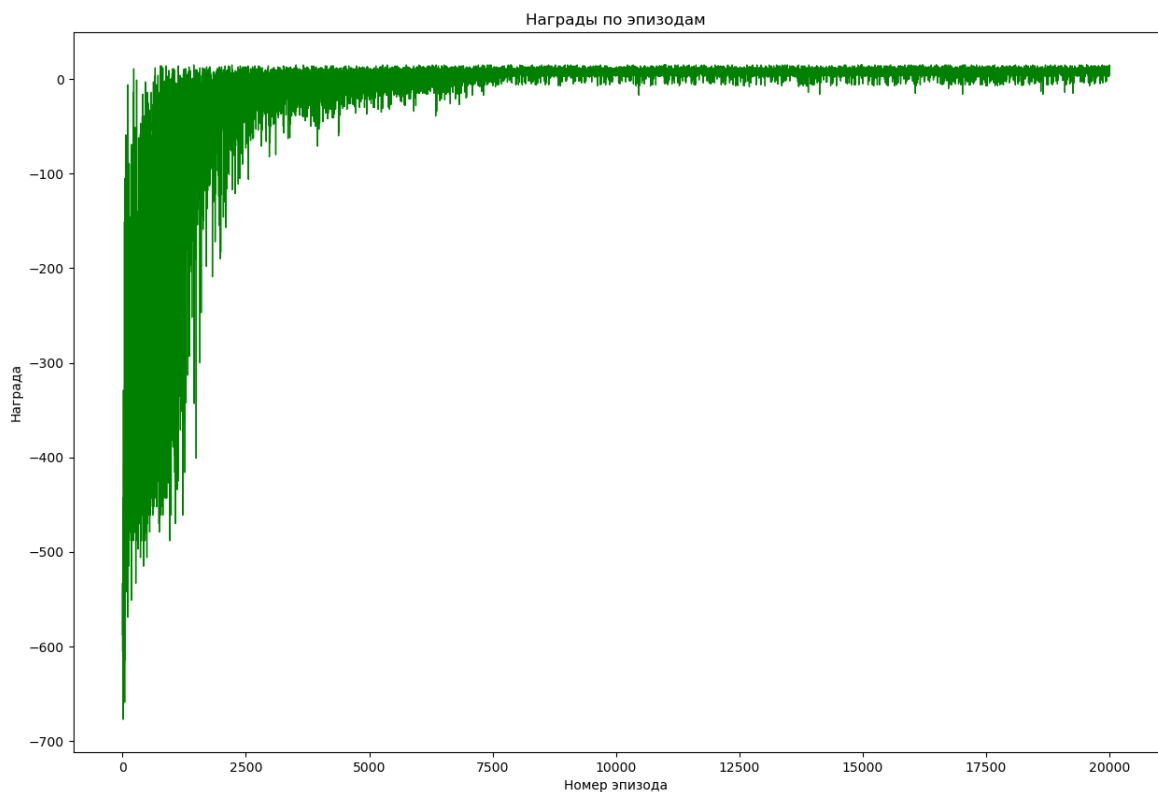


Двойное Q-обучение: $\text{eps}=0.4$, $\text{lr}=0.1$, $\text{gamma}=0.98$, $\text{num_episodes}=20000$

```

1 Вывод Q-матриц для алгоритма Двойное Q-обучение
2 Q1
3 [[ 0.          0.          0.          0.          0.          0.          ]
4  [-0.72267154  1.85252442  0.19552651  2.85144428  8.36234335 -6.59195419]
5  [ 7.18532871  6.59356075  5.37962679  5.5787384   13.27445578 -4.09009116]
6  ...
7  [-0.66467311 12.85534519 -1.95425954 -2.31569204 -5.48922472 -4.37459481]
8  [-3.22093218 -3.91954695 -3.58949805  1.68910615 -7.61442506 -8.67036163]
9  [ 3.79014035  3.88911184  5.13464853 18.25507162  0.68488221  1.74249151]]
10 Q2
11 [[ 0.          0.          0.          0.          0.          0.          ]
12  [-0.77584626  3.15691246  0.69250749  2.73240034  8.36234335 -5.35095443]
13  [ 4.20528463  7.68270196  4.62451724  3.56948691 13.27445578 -1.14602156]
14  ...
15  [-1.54713418 13.67037706  0.3485602  -0.64388233 -7.43964304 -4.11504528]
16  [-4.99805617 -4.01629887 -3.96321961  2.922531  -4.97259003 -3.3258903 ]
17  [ 2.08418592  0.94533048  1.13723617 18.48774274 -1.7645767  -0.50591363]]

```



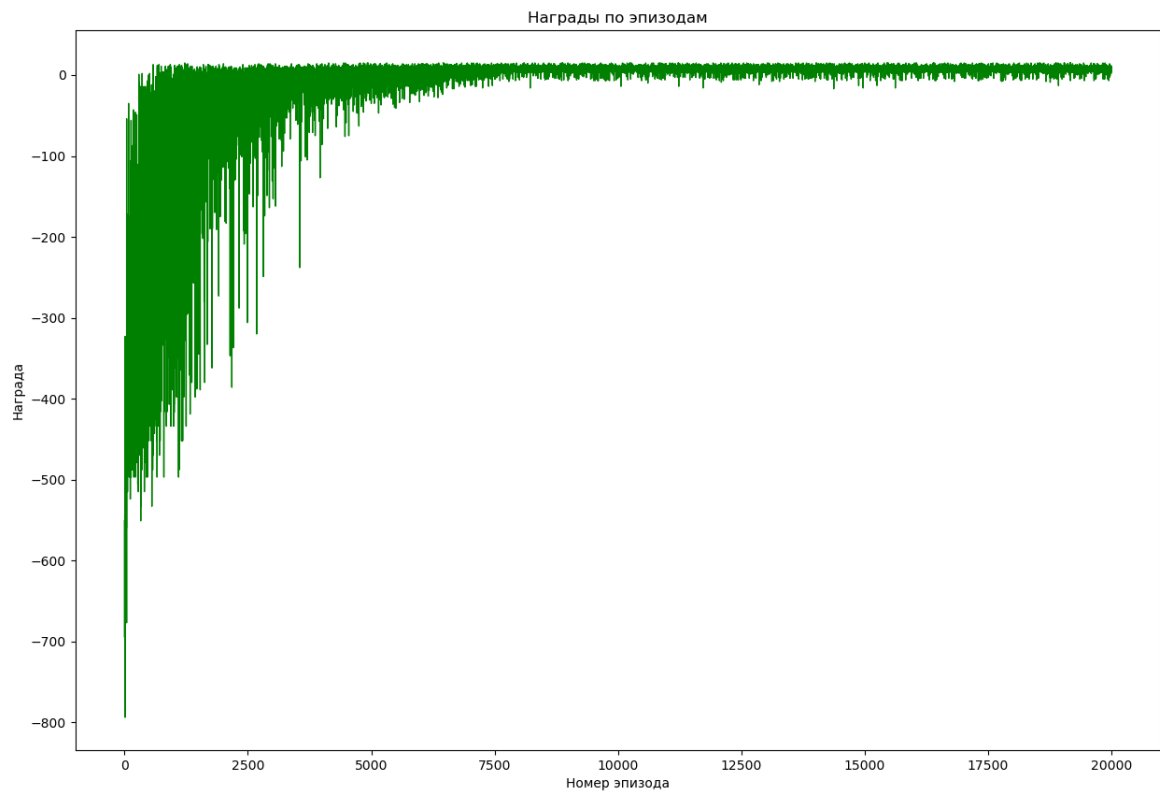
SARSA: $\text{eps}=0.4$, $\text{lr}=0.1$, $\text{gamma}=0.98$, $\text{num_episodes}=20000$

```

1 Вывод Q-матрицы для алгоритма SARSA
2 [[ 0.          0.          0.          0.          0.
3  0.          ]
4  [-8.66665356 -5.19580592 -3.60445722 -2.81367279  7.94897997
5  -10.02695965]
6  [-0.13055387 -0.3346179  0.34120444  5.34029633 13.19676101
7  -3.77660174]
8  ...

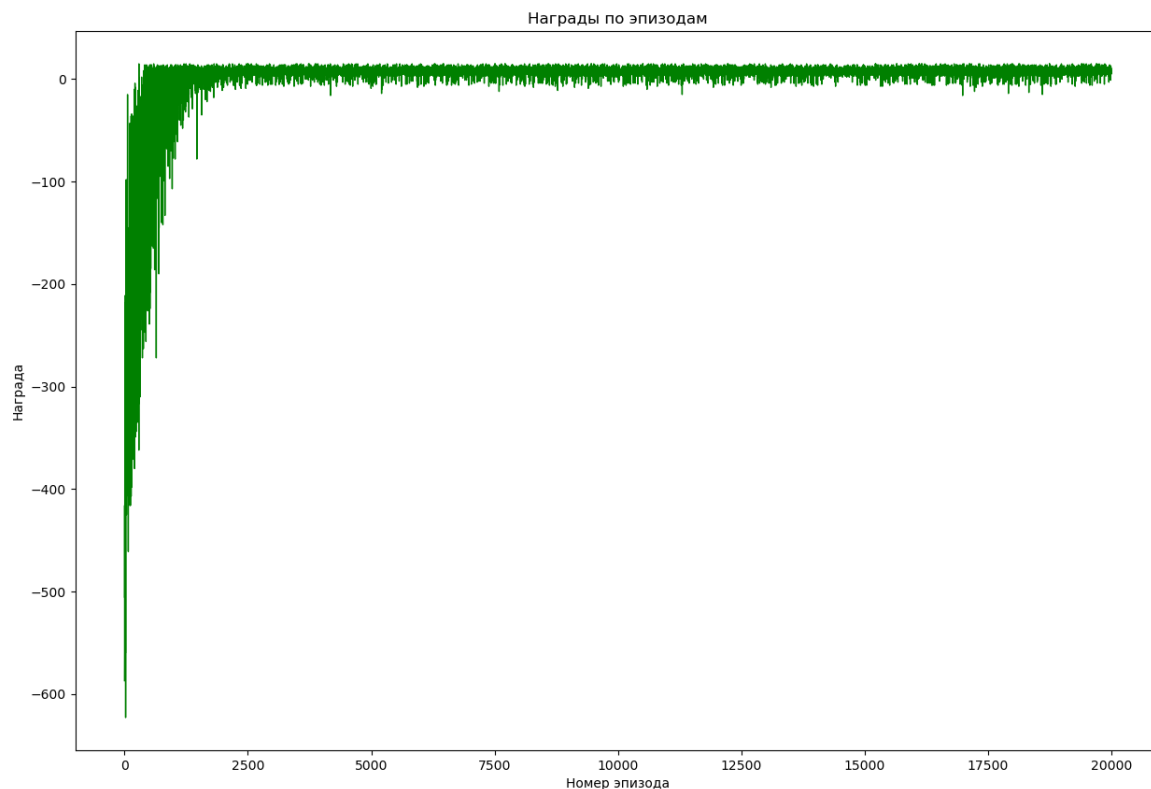
```

```
9  [ -0.65210655  10.87574541  -2.27322581  -2.06382686  -7.35227183
10   -4.84737165]
11  [ -8.37712909  -4.26957472  -8.15637214  -7.89962417  -14.04817395
12   -13.32447119]
13  [  8.81479907   7.35859204   6.83248527  18.47612868  -0.1078542
14   0.98660318]]
```



4. Вывод

На параметрах по умолчанию быстрее всего сходится Q-обучение. Gamma подобрана удачно, крупные изменения значения делают результат хуже. Быстрая сходимось получается на $\text{eps} = 0.1$:



SARSA также быстрее сходится при уменьшении eps . Learning rate и gamma подобраны удачно.