

МОСКОВСКИЙ ГОСУДАРСТВЕННЫЙ ТЕХНИЧЕСКИЙ УНИВЕРСИТЕТ
им. Н.Э. Баумана

Факультет «Информатика и системы управления»
Кафедра «Систем обработки информации и управления»

ОТЧЕТ

Лабораторная работа № 5
по дисциплине «Методы машинного обучения»

Тема: «Обучение на основе временны'х различий»

ИСПОЛНИТЕЛЬ:

группа ИУ5-25

Алексеев А С
ФИО

подпись

"__" _____ 2024 г.

ПРЕПОДАВАТЕЛЬ:

Гапанюк Ю Е
ФИО

подпись

"__" _____ 2024 г.

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).

SARSA очень напоминает Q-learning. Ключевое отличие SARSA от Q-learning заключается в том, что это алгоритм с политикой (on-policy). Это означает, что SARSA оценивает значения Q на основе действий, выполняемых текущей политикой, а не жадной политикой. Двойное обучение Q (Double Q-learning) представляет собой модификацию алгоритма обучения Q, который использует две функции ценности для уменьшения переоценки ценности действий. Этот метод помогает уменьшить переоценку ценности действий, которая может возникнуть в обычном алгоритме обучения Q.

```

In [1]: 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

# ***** 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())
        # Фла г т а тного завершения эпизода
        done = False
        # Флаг неш т а тного завершения эпизода
        truncated = False
        # Суммарная награда по эпизоду
        tot_rew = 0

        # По мере заполнения Q-м а т р и ц ы уменьшаем вероя т н о с т ь случайного выбора дейс т в и я
        if self.eps > self.eps_threshold:
            self.eps -= self.eps_decay

        # Выбор дейс т в и я
        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 дейс т в и е
            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] - self.Q[state][action])

            # Следующее состояние счи т а е м текущим
            state = next_state
            action = next_action
            # Суммарная награда за эпизод
            tot_rew += rew
            if (done or truncated):
                self.episodes_reward.append(tot_rew)

# ***** 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())
        # Флаги 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)

            # 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] - self.Q[state][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[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)

# ***** 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 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)

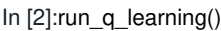
            # Выполняем шаг в среде
            next_state, rew, done, truncated, _ = self.env.step(action)

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

```


Награды по эпизодам



```
[ -12.65068321 -12.31234588 -12.31286591 -12.77198199]
[ -12.15372748 -11.54873869 -11.5487351 -12.81217421]
[ -11.44006738 -10.76414741 -10.76414828 -12.18884749]
[ -10.70446606 -9.96342994 -9.96342987 -11.40312183]
[ -9.90661774 -9.14635925 -9.14635928 -10.73306004]
[ -9.12623288 -8.31261184 -8.31261184 -9.9453808 ]
[ -8.29377416 -7.46184886 -7.46184886 -9.11819485]
[ -7.45824465 -6.59372334 -6.59372334 -8.29251683]
[ -6.59061771 -5.70388096 -5.70788096 -6.44546287]
[ -5.70257285 -4.807396016 -4.80396016 -7.55577764]
5 4.769214272 0.6814562 0.6814562 5.007200281
```


Награды по эпизодам



Вывод Q-матриц для алгоритма Двойное Q-обучение Q1

[illegible]

Q2

[[-16.18037032 -12.32854 -14.79330728 -14.75692416]
 [-13.10248251 -12.78202981 -11.54888163 -14.71939397]
 [-14.13280958 -13.63342026 -10.77521995 -13.2930076]
 [-13.82145874 -14.17277543 -10.0151928 -14.72737419]
 [-12.65989991 -12.91510771 -9.31383784 -13.68900962]
 [-11.78141116 -8.37774551 -11.49657163 -12.92774846]
 [-9.44792825 -7.46563286 -9.48989584 -10.64859996]
 [-8.86427966 -6.59374369 -7.82975565 -9.58021674]
 [-6.85613128 -6.95011406 -5.70788132 -8.13252839]
 [-7.63828248 -4.83430899 -5.1491614 -7.04143969]
 [-5.25744611 -3.88749159 -7.31534274 -7.32988921]
 [-4.39425883 -4.38959936 -2.94040938 -5.2512676]
 [-13.29368973 -11.54888054 -11.63499601 -12.40434605]
 [-12.31933251 -10.76424416 -10.76416381 -12.31794733]
 [-11.71038621 -9.99899947 -9.96343246 -11.57219514]
 [-11.03629509 -9.23328926 -9.14635966 -10.77847981]
 [-10.69183139 -8.31261189 -8.55205744 -10.22938403]
 [-9.77498619 -7.47737233 -7.46184887 -9.15648549]
 [-8.75236126 -6.59372334 -7.05389519 -8.74021723]
 [-7.57034127 -5.70788096 -5.72942051 -7.5513079]
 [-6.59491342 -4.80914676 -4.80396016 -6.59707052]
 [-6.72107989 -6.65139743 -3.881592 -5.98369779]
 [-4.91406581 -2.9404 -2.91675167 -4.78424944]
 [-3.89050651 -2.94610952 -1.98 -6.01593777]
 [-12.31790293 -10.76416381 -12.31790293 -11.54888054]
 [-11.54888054 -9.96343246 -11.31790293 -11.54888054]
 [-10.76416381 -9.14635966 -11.31790293 -10.76416381]

