

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 в автоматизированных
системах»

Тема: «Обучение на основе временны'х различий»

ИСПОЛНИТЕЛЬ:

группа ИУ5-22М

Калюта Н.И.
ФИО

подпись

"30" 05 2024 г.

ПРЕПОДАВАТЕЛЬ:

ФИО

подпись

" " _____ 2024 г.

Москва - 2024

Цель лабораторной работы:

Ознакомление с базовыми методами обучения с подкреплением на основе временных различий.

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).

Установка необходимых библиотек и зависимостей для работы с Atari-играми

```
!rm -rf /dev/null 2>&1" to see what is going on under the hood
!apt-get update > /dev/null 2>&1
!pip install gym pyvirtualdisplay > /dev/null 2>&1
!apt-get install -y xvfb python-opengl ffmpeg > /dev/null 2>&1
!apt-get install cmake > /dev/null 2>&1
!pip install --upgrade setuptools 2>&1
!pip install ez_setup > /dev/null 2>&1

!apt-get update > /dev/null 2>&1
!apt-get install cmake > /dev/null 2>&1
!pip install --upgrade setuptools 2>&1
!pip install ez_setup > /dev/null 2>&1

!pip install swig==4.1.1
!pip install Box2D==2.3.2
!pip install box2d-kengz==2.3.3
!pip install pygame==2.2.0
!pip install ale_py==0.8.1
!pip install pygamelet==1.5.11

!pip install -U colabgymrender
!pip install imageio==2.4.1
!pip install --upgrade AutoROM
!AutoROM --accept-license
!pip install gymnasium[atari]==0.28.1

!wget http://www.atarimania.com/roms/Roms.rar
!unrar x -o+ /content/Roms.rar >/dev/null
!python -m atari_py.import_roms /content/ROMS >/dev/null

!pip install pyvirtualdisplay > /dev/null 2>&1

!pip install --upgrade gym
```

Установка библиотек, необходимые для работы с графикой и видео

```
[ ] !apt-get install xvfb
!apt-get install python3-opengl ffmpeg
```

Создать папку с названием "video"

```
[ ] !mkdir -p video
```

Импорт библиотек

```
[ ] import gymnasium as gym
from gymnasium import logger as gymlogger
from gymnasium.wrappers.record_video import RecordVideo

gymlogger.set_level(40) #error only
import tensorflow as tf
import numpy as np
import random
import matplotlib
import matplotlib.pyplot as plt
%matplotlib inline
from pprint import pprint
from tqdm import tqdm
import math
import uuid
import glob
import io
import base64
from IPython.display import HTML

from IPython import display as ipythondisplay
```

```
[ ] gym.__version__
```

```
↗ '0.28.1'
```

▼ Дополнительные функции

Функции для работы с видеозаписями, создает среду Gym с возможностью записи видео и выводит информацию о среде

```
def show_video(folder_name):
    mp4list = glob.glob(f'{folder_name}/*.mp4')
    if len(mp4list) > 0:
        mp4 = mp4list[0]
        video = io.open(mp4, 'r+b').read()
        encoded = base64.b64encode(video)
        ipythondisplay.display(HTML(data='{}'.format(encoded.decode('ascii'))))
    else:
        print("Could not find video")

def wrap_env(env, folder_name):
    env = RecordVideo(env, folder_name, step_trigger = lambda episode_number: True)
    return env

def create_environment(name):
    folder_name = f"./video/{name}/{uuid.uuid4()}"
    env = wrap_env(gym.make(name, render_mode="rgb_array"), folder_name)
    spec = gym.spec(name)
    print(f"Action Space: {env.action_space}")
    print(f"Observation Space: {env.observation_space}")
    print(f"Max Episode Steps: {spec.max_episode_steps}")
    print(f"Nondeterministic: {spec.nondeterministic}")
    print(f"Reward Range: {env.reward_range}")
    print(f"Reward Threshold: {spec.reward_threshold}")
    return env, folder_name
```

Создает невидимый виртуальный дисплей размером 1400x900 пикселей

```
[ ] from pyvirtualdisplay import Display
display = Display(visible=0, size=(1400, 900))
display.start()
```

```
↗ <pyvirtualdisplay.display.Display at 0x7ec27399b430>
```

▼ Задание агентов

Реализация базового агента

```
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

```

Реализация алгоритма обучения с подкреплением SARSA (State-Action-Reward-State-Action) для агента

```

[ ] 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 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, info = 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] - self.Q[state][action])

# 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)

```

Реализация алгоритма обучения с подкреплением Q-Learning

```

[ ] 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 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, info = self.env.step(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)

```

Реализация алгоритма обучения с подкреплением "Double Q-Learning". Алгоритм Double Q-Learning является вариантом алгоритма Q-Learning, который помогает уменьшить переоценку значений Q и улучшить стабильность обучения.

```
[ ] 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=0.0005
        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)

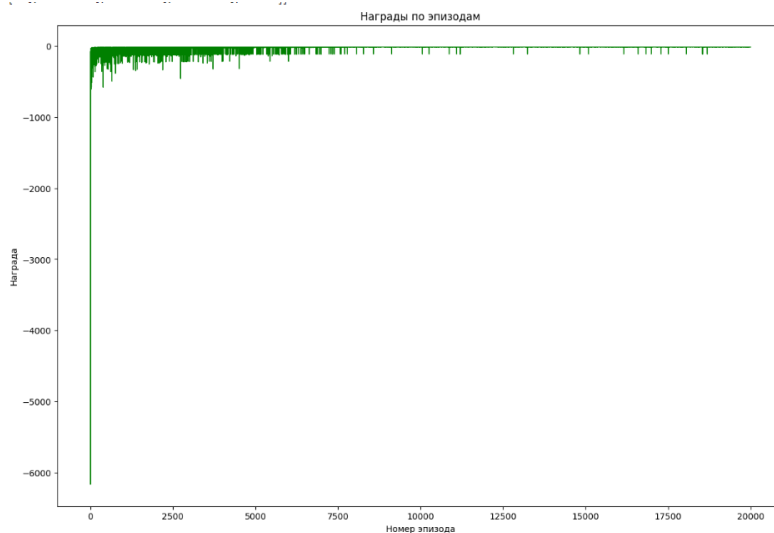
                # Выполняем шаг в среде
                next_state, rew, done, truncated, _ = self.env.step(action)

                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)
```

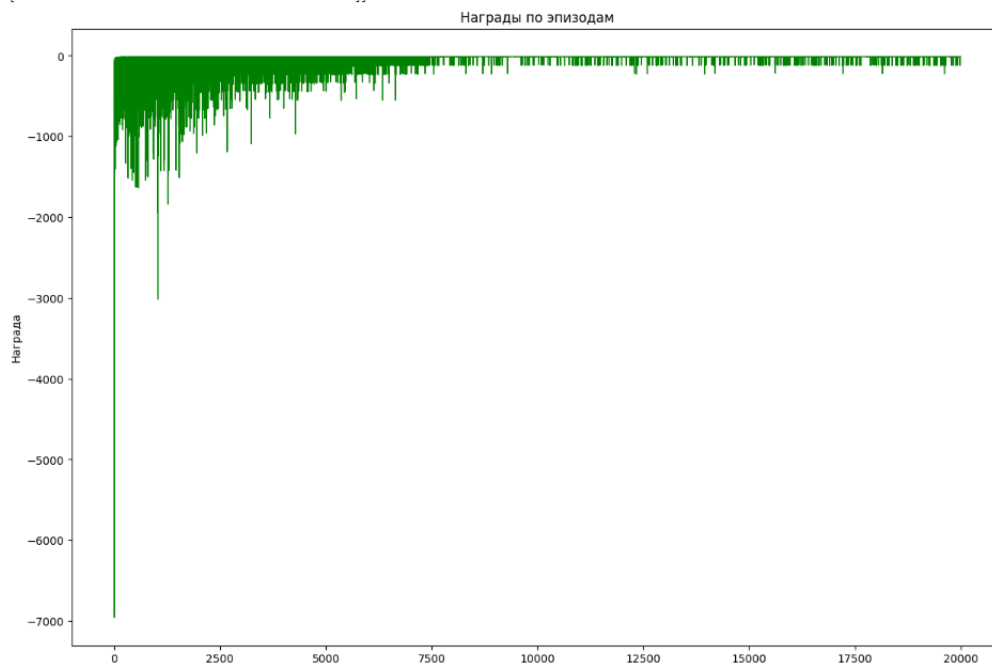
```
[ ] def play_agent(agent):
    """
    Проигрывание сессии для обученного агента
    """
    env2, folder = create_environment('CliffWalking-v0')
    state = env2.reset()[0]
    done = False
    while not done:
        action = agent.greedy(state)
        next_state, reward, terminated, truncated, info = env2.step(action)
        env2.render()
        state = next_state
        if terminated or truncated:
            done = True
    show_video(folder)
    env2.close()
```

Реализация обучения и тестирования агента, использующего алгоритм SARSA, в среде CliffWalking-v0.

[illegible]

Реализует полный цикл обучения и тестирования агента с использованием алгоритма Q-Learning

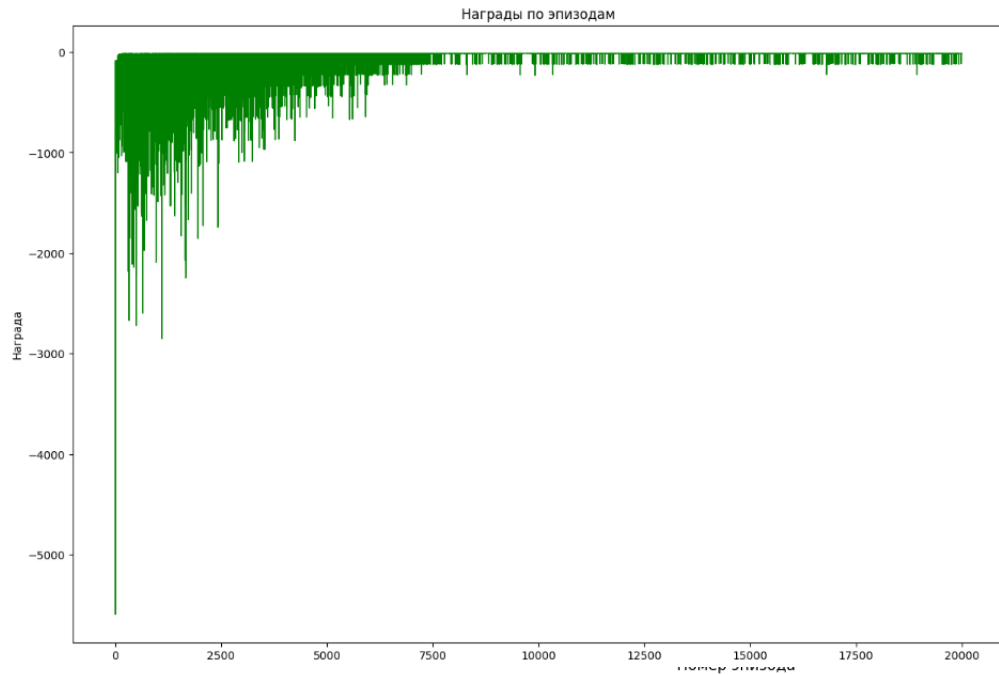
```
[ ] env = gym.make('CliffWalking-v0')
    agent = QLearning_Agent(env)
    agent.learn()
    agent.print_q()
    agent.draw_episodes_reward()
    play_agent(agent)
```

[illegible]

Реализует полный цикл обучения и тестирования агента с использованием алгоритма Double Q-Learning

```
[ ] env = gym.make('CliffWalking-v0')
agent = DoubleQLearning_Agent(env)
agent.learn()
agent.print_q()
agent.draw_episodes_reward()
play_agent(agent)
```

[illegible]



```
Action Space: Discrete(4)
Observation Space: Discrete(48)
Max Episode Steps: None
Nondeterministic: False
Reward Range: (-inf, inf)
Reward Threshold: None
Moviepy - Building video /content/video/CliffWalking-v0/fd24cdc1-bde8-4623-8962-e879695b0be2/r1-video-step-0.mp4.
Moviepy - Writing video /content/video/CliffWalking-v0/fd24cdc1-bde8-4623-8962-e879695b0be2/r1-video-step-0.mp4

Moviepy - Done !
Moviepy - video ready /content/video/CliffWalking-v0/fd24cdc1-bde8-4623-8962-e879695b0be2/r1-video-step-0.mp4
```

Вывод

В ходе выполнения лабораторной работы ознакомился с тремя алгоритмами обучения с подкреплением: SARSA, Q-Learning и Double Q-Learning.

Реализовал эти алгоритмы в коде и провели их сравнение на задаче "обрыва" (CliffWalking-v0). Наблюдал, как агенты, обученные с помощью каждого из алгоритмов, справляются с задачей. Сравнил эффективность обучения, изучив графики накопленной награды за каждый эпизод для каждого алгоритма. В итоге получил практическое представление о работе алгоритмов обучения с подкреплением, их отличиях и о том, как они справляются с задачей "обрыва".