

МОСКОВСКИЙ ГОСУДАРСТВЕННЫЙ ТЕХНИЧЕСКИЙ УНИВЕРСИТЕТ
им. Н.Э. Баумана

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

ОТЧЕТ

Лабораторная работа № 4
по дисциплине «Методы машинного обучения в АСОИУ»

ИСПОЛНИТЕЛЬ:

группа ИУ5-22М

Воронцова А.В.

ФИО

подпись

"__" "__" 2024 г.

ПРЕПОДАВАТЕЛЬ:

Балашов А. М.

ФИО

подпись

"__" "__" 2024 г.

Москва - 2024

Задание

На основе рассмотренного на лекции примера реализуйте алгоритм Policy Iteration для любой среды обучения с подкреплением (кроме рассмотренной на лекции среды Toy Text / Frozen Lake) из библиотеки Gym (или аналогичной библиотеки).

```
pip install gym

Requirement already satisfied: gym in /usr/local/lib/python3.10/dist-packages (0.25.2)
Requirement already satisfied: numpy>=1.18.0 in /usr/local/lib/python3.10/dist-packages (from gym) (1.25.2)
Requirement already satisfied: cloudpickle>=1.2.0 in /usr/local/lib/python3.10/dist-packages (from gym) (2.2.1)
Requirement already satisfied: gym-notices>=0.0.4 in /usr/local/lib/python3.10/dist-packages (from gym) (0.0.8)

import numpy as np
import gym
import matplotlib.pyplot as plt

def policy_evaluation(policy, env, discount_factor=0.95, theta=1e-5):
    V = np.zeros(env.observation_space.n)
    while True:
        delta = 0
        for s in range(env.observation_space.n):
            v = 0
            for a, action_prob in enumerate(policy[s]):
                for prob, next_state, reward, done in env.P[s][a]:
                    v += action_prob * prob * (reward + discount_factor * V[next_state])
            delta = max(delta, abs(V[s] - v))
            V[s] = v
        if delta < theta:
            break
    return V

def policy_improvement(V, env, discount_factor=0.95):
    policy = np.zeros([env.observation_space.n, env.action_space.n])
    for s in range(env.observation_space.n):
        q = np.zeros(env.action_space.n)
        for a in range(env.action_space.n):
            for prob, next_state, reward, done in env.P[s][a]:
                q[a] += prob * (reward + discount_factor * V[next_state])
        best_action = np.argmax(q)
        policy[s, best_action] = 1.0
    return policy

def policy_iteration(env, discount_factor=0.95):
```

[illegible]

[illegible]