

```

from cartpole import CartPoleEnv
import numpy as np

import random

env = CartPoleEnv()

def discretize(val,bounds,n_states):
    discrete_val = 0
    if val <= bounds[0]:
        discrete_val = 0
    elif val >= bounds[1]:
        discrete_val = n_states-1
    else:
        discrete_val = int(round((n_states-1)*((val-bounds[0])/(bounds[1]-bounds[0]))))
    return discrete_val

def discretize_state(vals,s_bounds,n_s):
    discrete_vals = []
    for i in range(len(n_s)):
        discrete_vals.append(discretize(vals[i],s_bounds[i],n_s[i]))
    return np.array(discrete_vals,dtype=np.int)

#parametry dyskretyzacji
n_s = np.array([7,7,7,7])
n_a = env.action_space.n

s_bounds = np.array(list(zip(env.observation_space.low,env.observation_space.high)))
s_bounds[1] = (-1.0,1.0)
s_bounds[3] = (-1.0,1.0)

s_bounds = np.dtype('float64').type(s_bounds)

Q = np.zeros(np.append(n_s,n_a))

print(Q.shape)

def epsilon_greedy_action_from_Q(env, state, epsilon=0.2):
    if np.random.random() < epsilon:
        action = env.action_space.sample()
    else:
        action = np.argmax(Q[tuple(state)])
    return action

def SARSA_Q(env, episodes=1000, gamma=0.9, alpha=0.3):

    #tablica do której zapisujemy sumę nagród z każdego epizodu
    Rewards = []

    for i in range(episodes):

        if (i%100)==0:
            print("episode=" + i)

```

```

print('SARSA Q-learning')

obs = env.reset()
S = discretize_state(obs,s_bounds,n_s)

episode_reward = 0
finished = False

a = epsilon_greedy_action_from_Q(env,S)

time_step=0

#zakończenie epizodu gdy 'finished == True' lub 'ilość kroków == 200'
while not finished and not time_step==200:

    #DO UZUPEŁNIENIA
    obs, reward, finished, info = env.step(a)
    next_S = discretize_state(obs,s_bounds,n_s)

    next_A =epsilon_greedy_action_from_Q(env,S)

    Q[(S,a)] = Q[(S,a)] + alpha * (reward + gamma * Q[(next_S, next_A)] - Q[(S,a)])

    #DO UZUPEŁNIENIA

    S = next_S
    a = next_A

    #sumujemy wszystkie nagrody zdobyte w danym epizodzie
    episode_reward += reward

    Rewards.append(episode_reward)

return Q, Rewards

learning_episodes = 3000

_,R = SARSA_Q(env,learning_episodes)

#Wyliczamy średnią nagrodę - ilość epizodów (learning_episodes) dzielimy na 100
#i uśredniamy nagrody z kolejnych (learning_episodes/100) epizodów.
meanR= []
for i in range(100):
    meanR.append(np.mean(R[int(learning_episodes/100)*i:int(learning_episodes/100)*(i+1)]))

#wykres pokazuje nagrody zdobyte w 100 kolejnych epizodach wybranych z wszystkich epizodów
#oraz jak zmieniała się średnia nagroda zdobyta przez agenta

import matplotlib.pyplot as plt
x_data = range(0,100)
plt.plot(x_data,R[:int(learning_episodes/100)],label="reward")
plt.plot(x_data,meanR,label="mean reward")
plt.title('CartPole: SARSA')
plt.xlabel('Episode')

```

```
plt.ylabel('Reward')
plt.legend()
plt.show()
```

```

/usr/local/lib/python3.6/dist-packages/gym/logger.py:30: UserWarning: WARN: Box bound
  warnings.warn(colorize('%s: %s'%( 'WARN', msg % args), 'yellow'))
(7, 7, 7, 7, 2)
episode= 0
episode= 100
episode= 200
episode= 300
episode= 400
episode= 500
episode= 600
episode= 700
episode= 800
episode= 900
episode= 1000
episode= 1100
episode= 1200
episode= 1300
episode= 1400
episode= 1500
episode= 1600
episode= 1700
episode= 1800
episode= 1900
episode= 2000
episode= 2100
episode= 2200
episode= 2300
episode= 2400
episode= 2500
episode= 2600
episode= 2700
episode= 2800
episode= 2900

```



