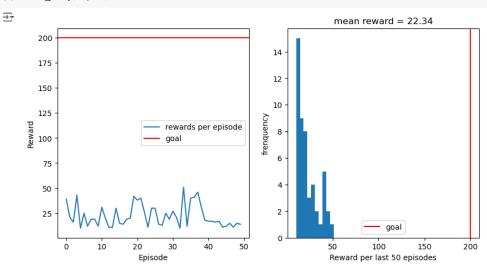


```
plt.pause(0.02)
        display.clear_output(wait=True)
        fig, ax = plt.subplots(1,2, figsize=(10,5)) # 2 subplots in one figure
        fig.suptitle(title)
        # plot 1 — rewards per episode
        ax[0].plot(values, label="rewards per episode")
        ax[0].axhline(200, c='red', label='goal')
ax[0].set_xlabel('Episode')
        ax[0].set_ylabel('Reward')
        ax[0].legend()
        # plot 2 — mean reward of the 1st 50 episodes
        ax[1].set_title("mean reward = {}".format(sum(values[-50:])/50))
        ax[1].hist(values[-50:])
        ax[1].axvline(200, c='red', label='goal')
        ax[1].set_xlabel('Reward per last 50 episodes')
ax[1].set_ylabel('frenquency')
        ax[1].legend()
plt.show()
```

[] random_policy(env, 50)



✓ 2. Deep Q Learning

We defined a deep Q learning class which contains the structure of a neural network, a update function and a predict function.

```
[8] print(torch._path_)
```

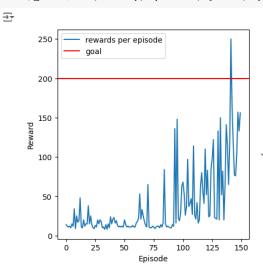
['/usr/local/lib/python3.10/dist-packages/torch']

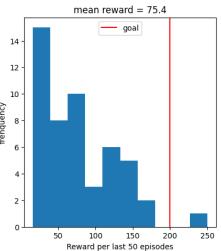
We defined a **DQN** function with the epsilon greedy policy.

and should_run_async(code)

```
[18] def DQN_process(env, model, episodes, gamma=0.9, epsilon=0.3, epsilon_decay=0.99, title="DQN_model"):
            gamma: Markov procession gamma
            epsilon: we use epsilon greedy algorithm
            epsilon_decay: epsilon decay rate
            rewards = []
            for _ in range(episodes):
                   state = env.reset()
                    total reward = 0
                   done = False
                    while not done:
                           # epsilon greedy: if the random value less than epsilon, choose the random value, otherwise use the predict value
                           if random.random() < epsilon:</pre>
                                  action = env.action_space.sample()
                                  q_values = model.predict(state)
                                  action = torch.argmax(q_values).item() # find the maximum item in q_values by their index
                           next_state, reward, done, info = env.step(action) # the four parameters by doing step"action"
                           total_reward += reward
                           q_values = model.predict(state)
                           # q value of the current step = former reward + gamma * the best prediction
# tell the agend the episode if it is end
                           q_values[action] = (reward + (0 if done else gamma * torch.max(model.predict(next_state)).item()) )
                           model.update(state, q_values) # update the model base on the new q value
                           state = next_state # jump to the next state
                    epsilon = max(epsilon * epsilon_decay, 0.001) # epsilon decay until 0.001
                    rewards.append(total_reward)
                    plot_result(rewards)
```

// [17] dqn = DQN(state_dim=4, action_dim=2, hidden_dim = 64, lr=0.001)
DQN_process(env=env, model=dqn, episodes=150, gamma=0.9, epsilon=0.5)





- [] 开始借助 AI 编写或<u>生成</u>代码。
- [] 开始借助 AI 编写或生<u>成</u>代码。

Colab 付费产品 - 在此处取消合同

✓ 已连接到"Python 3 Google Compute Engine 后端"

×