

بسم الله الرحمن الرحيم

تمرین یادگیری تقویتی درس کنترل هوشمند در مکاترونیک

اعضای گروه:

محمدامین محتشمی پور، مرتضی جداری، پویان سلیمانی، محمد میرزایی، وانیلا ملک زاهدی

نام استاد:

دکتر ایمان شریفی

بهار ۱۴۰۳



دانشگاه صنعتی امیرکبیر
(پلی تکنیک تهران)

فهرست

٤	١	سوالات تئوری
٤	١/١	
٤	١/٢	
٥	١/٣	
٨	٢	پروژه عملی
٨	٢/١	مقدمه
٩	٢/٢	ساخت محیط
١٠	٢/٣	یادگیری محیط
١٢	٢/٤	رسم مدل

۱ سوالات تئوری

۱/۱

Sample Efficiency یا همان بهینگی در نمونه گیری به الگوریتمی گفته می شود که از هر نمونه بیشترین بهره را ببرد. یعنی بتواند از هر قطعه از تجربه ای که تولید می شود استفاده بهینه ببرد و سریعاً سیاست خود را بهبود ببخشد.

برخلاف آن، اگر یک الگوریتم از نمونه هایی که دارد تجربیات درستی به دست نیاورد، آن الگوریتم بهینه نخواهد بود.

یکی از روش هایی که برای این امر استفاده شده است، نمونه برداری اهمیتی یا Importance sampling می باشد که یک روش مونته کارلو برای ارزیابی ویژگی های یک توزیع خاص است.

همچنین از دیگر راه های افزایش Sample Efficiency استفاده از رویکرد مبتنی بر Meta-Learning می باشد که هدف آن این است که مدل بتواند با داده های جدید سریع تر سازگار شود و نیاز به نمونه های کمتری داشته باشد. [۱]

به علاوه یکی دیگر از رویکردهای افزایش Sample Efficiency استفاده از Auxiliary Tasks یا همان وظایف کمکی می باشد. [۱]

همچنین یکی دیگر از راه های افزایش Sample Efficiency استفاده از بافر تجارت یا Replay Buffer می باشد.

۱/۲

بافر تجارت (Replay Buffer) یک تکنیک مهم در یادگیری تقویتی است که به افزایش sample efficiency کمک می کند.

در یادگیری تقویتی، عامل (agent) با تعامل با محیط، تجربیات (transitions) مختلفی را به صورت توالی دنباله ای (sequence) کسب می کند. این تجربیات شامل وضعیت فعلی، اقدام انجام شده، پاداش دریافتی و وضعیت بعدی هستند.

بافر تجارت یک حافظه است که این تجربیات را ذخیره می کند. در طول آموزش، عامل به جای استفاده از تجربیات جدید، از این بافر به صورت تصادفی نمونه برداری می کند و از آنها برای به روزرسانی شبکه عصبی استفاده می کند.

این رویکرد به چند دلیل باعث افزایش sample efficiency می شود:

۱. تکرار تجربیات: استفاده مکرر از تجربیات ذخیره شده در بافر به جای تجربیات جدید، به عامل کمک می کند تا از آنها بیشترین استفاده را ببرد.

۲. شکستن همبستگی: در یادگیری تقویتی، توالی تجربیات ممکن است همبستگی داشته باشند. بافر تجارت با نمونه برداری تصادفی از تجربیات، این همبستگی را می شکند و به پایداری آموزش کمک می کند.

۳. تنوع تجربیات: بافر تجارت می تواند تجربیات متنوعی را ذخیره کند. این به عامل کمک می کند تا مدل خود را به طور کلی تری آموزش دهد.

در مجموع، بافر تجارت یک تکنیک قدرتمند برای افزایش sample efficiency در یادگیری تقویتی است که با ذخیره و استفاده مجدد از تجربیات، به عامل کمک می کند تا با داده های کمتری بهتر یاد بگیرد.

در این مقاله از دو الگوریتم (Trust Region Policy Optimization (TRPO و Deep Q-Network with NAF (DQN استفاده شده و نتایج آن گزارش شده است و هدف آن مقایسه‌ی این دو الگوریتم با دو الگوریتم دیگر یعنی Deep Deterministic Policy Gradient (DDPG و Vanilla Policy Gradient (VPG می‌باشد. سیاست این الگوریتم به‌صورت زیر فرموله شده است و با شبکه عصبی عمیق یارمتر شده است.

حالت های این الگوریتم به صورت زیر می باشد:

$$q_{ur5} = [q_{\setminus} \ q_{\setminus} \ q_{\setminus} \ q_{\setminus} \ q_{\setminus}]^T, q_{grip} = [q_{\setminus} \ ... \ q_{\setminus} \ q_{\setminus} \ ... \ q_{\setminus} \ q_{\setminus} \ ... \ q_{\setminus}]^T,$$

$$a_{ur5} = [m_{\setminus} \ m_{\setminus} \ m_{\setminus} \ m_{\setminus} \ m_{\setminus}]^T, a_{grip} = [m_{\setminus} \ m_{\setminus} \ m_{\setminus} \ m_{\setminus}]^T$$

که q ها، بردارهای موقعیت گیر و a ها بردارهای عمل (action) گیر هستند.

دو الگوریتم بیان شده در مقاله به صورت زیر هستند:

Algorithm 1 Trust Region Policy Optimization (TRPO)

```

1: Randomly initialize policy parameters  $\theta_0$ 
2: for  $k = 0, 1, 2, \dots$ , do
3:   Collect set of  $N_\tau$  trajectories under policy  $\pi_{\theta_k}$ 
4:   Estimate advantages with GAE( $\lambda$ ) and fit  $V^{\pi_{\theta_k}}$ 
5:   Estimate policy gradient


$$\hat{g}_k = \frac{1}{N_\tau} \sum_{i=1}^{N_\tau} \sum_{t=0}^{T-1} \gamma^t \nabla_{\theta} \log \pi_{\theta}(a_{t,i} | s_{t,i}) \Big|_{\theta=\theta_k} \hat{A}(s_{t,i}, a_{t,i})$$


6:   Estimate  $\hat{H}_k = \nabla_{\theta}^2 D_{KL}(\pi_{\theta} || \pi_{\theta_k}) \Big|_{\theta=\theta_k}$ 
7:   Compute  $\hat{H}_k^{-1} \hat{g}_k$  with CG algorithm
8:   Compute policy step  $\Delta_k = \sqrt{\frac{2\delta_D}{\hat{g}_k^T \hat{H}_k^{-1} \hat{g}_k}} \hat{H}_k^{-1} \hat{g}_k$ 
9:   for  $l = 1, 2, \dots, L$  do
10:    Compute candidate update  $\theta^c = \theta_k + \nu^l \Delta_k$ 
11:    if  $\mathcal{L}_{\pi_{\theta_k}}(\pi_{\theta^c}) \geq 0$  and  $D_{KL}(\pi_{\theta^c} || \pi_{\theta_k}) \leq \delta_D$ 
12:      then
13:        Accept candidate  $\theta_{k+1} = \theta^c$ 
14:      end if
15:    end for
16:  end for

```

Algorithm 2 Deep Q-Network with NAF (DQN-NAF)

```

1: Randomly initialize  $Q$  and target  $Q'$  with  $\theta^{Q'} \leftarrow \theta^Q$ 
2: Allocate Replay buffer  $\mathcal{R}$ 
3: for episode  $1, \dots, N_\tau$  do
4:   for  $t = 1, \dots, T$  do
5:     Execute  $a_t = \mu_{\theta^\mu}(s_t)$ 
6:     Store in  $\mathcal{R}$  transition  $(s_t, a_t, r_t, s_{t+1})$ 
7:   for iteration  $k = 1, \dots, K_Q$  do
8:     Sample minibatch of  $N_b$  transitions from  $\mathcal{R}$ 
9:     Set targets  $y_i = r_i + \gamma V_{\theta^{Q'}}(s_{t+1})$ 
10:    Update  $\theta^Q$  by minimizing loss


$$L(\theta^Q) = \frac{1}{N_b} \sum_{i=1}^{N_b} \left( y_i - Q_{\theta^Q}(s_i, a_i) \right)^2$$


11:    Update target network
12:     $\theta^{Q'} \leftarrow \xi \theta^Q + (1 - \xi) \theta^{Q'}$ 
13:  end for
14: end for

```

این مقاله الگوریتم ها را در دو وظیفه‌ی مختلف بررسی کرده است. یکی رسیدن به نقطه تصادفی و دیگری برداشتن و گذاشتن شیء (peak and place).

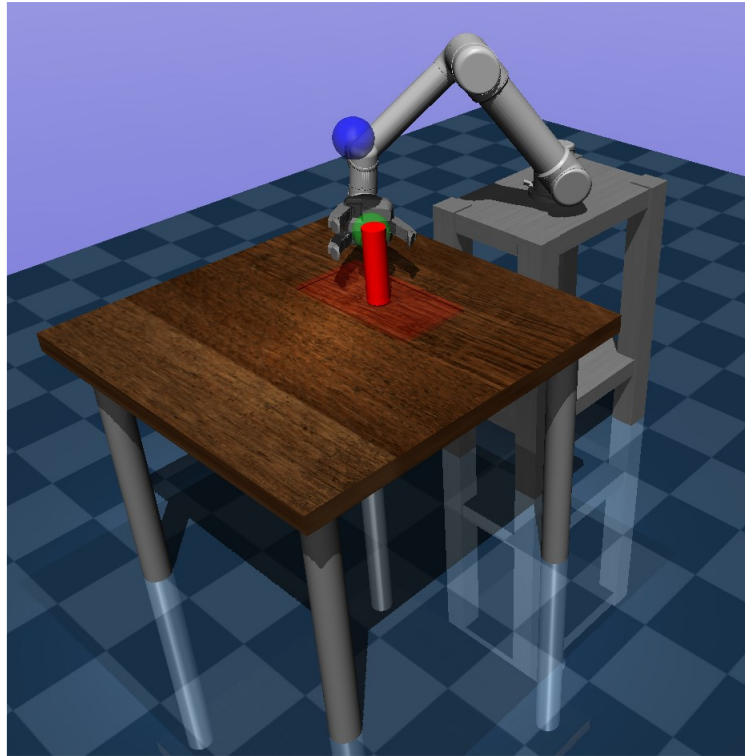
برای رسیدن به نقطه تصادفی، تابع جایزه به صورت زیر تعریف شده است:

$$r_R(s, a) = -||p_{\text{goal}} - p_{\text{ee}}|| - c_a ||a||$$

و همچنین برای برداشتن و گذاشتن شیء، تابع جایزه به شکل زیر است:

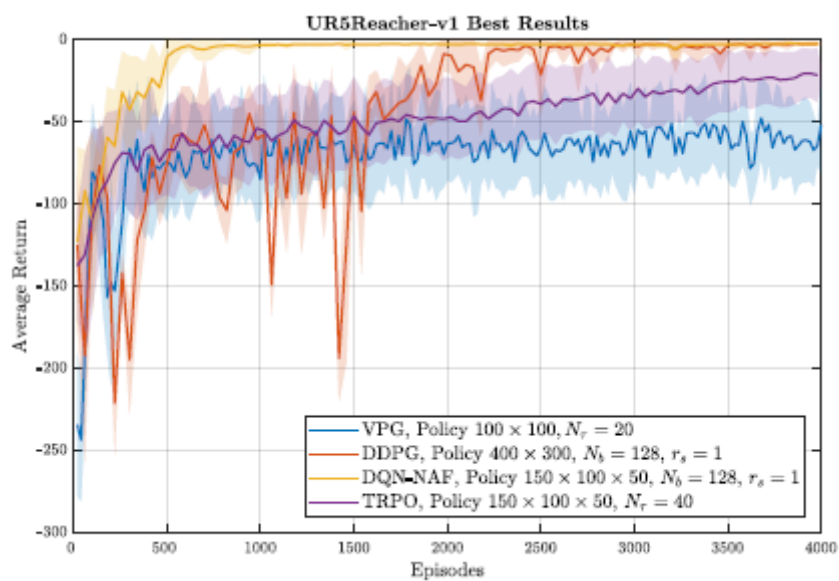
$$r_P(s, a) = \frac{1}{1 + \sum_{j=1}^3 c_j d_j} - c_a ||a||$$

محیط شبیه‌سازی به شکل زیر است:



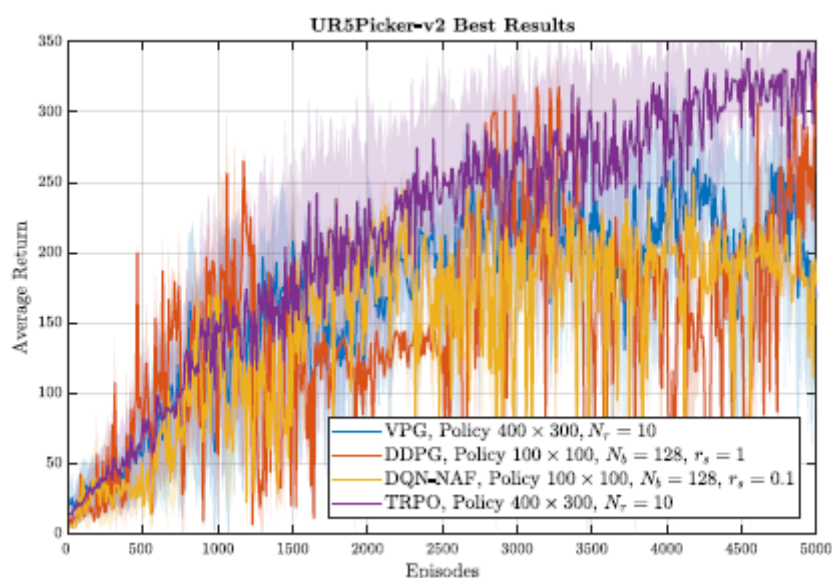
نتایج به دست آمده برای دو الگوریتم پیشنهاد شده و همچنین دو الگوریتم مورد مقایسه در حالت رسیدن به نقطه دلخواه به شکل زیر بیان شده است:

Algorithm	Episodes Req.	Final Avg Return	Max Return
VPG	260	-60.90 ± 24.70	-47.65 ± 24.71
DDPG	2860	-2.83 ± 2.57	-2.13 ± 1.21
DQN-NAF	700	-2.36 ± 1.23	-2.18 ± 1.04
TRPO	2360	-22.87 ± 16.27	-20.27 ± 14.83



همچنین این نتایج برای حالت برداشتن شیء به صورت زیر است:

Algorithm	Episodes Req.	Final Avg Return	Max Return
VPG	740	187.11 ± 70.63	273.95 ± 18.21
DDPG	4950	257.02 ± 12.59	321.24 ± 14.37
DQN-NAF	1980	173.08 ± 9.38	255.39 ± 5.73
TRPO	4980	324.97 ± 43.35	344.01 ± 17.74



[۱] Towards Sample Efficient Reinforcement Learning; Yang Yu

[۲] Robotic Arm Control and Task Training through Deep Reinforcement Learning; Andrea Franceschetti, Elisa Tosello, Nicola Castaman, and Stefano Ghidoni

۲ پروژه عملی

۲/۱ مقدمه

در این پروژه قرار است یک ربات سه درجه آزادی صفحه‌ای با استفاده از مدل یادگیری تقویتی، برای رسیدن به یک نقطه‌ی خاص آموزش ببیند.

۲/۲ ساخت محیط

برای ساخت محیط یادگیری، از کتابخانه‌ی `gymnasium` استفاده شده است. با کمک این کتابخانه، یک کلاس جدید به نام `ThreeLinkRobotEnv` که از سوپر کلاس `gym.Env` ارث بری شده است، ساخته شده است و توابع لازم آنها یعنی `__init__`، `reset`، `step` و `render` فراخوانی شده اند. در تابع `step` مقادیر مربوط به پاداش قرار داده شده است. پاداش این ربات با استفاده از تابع $\frac{20}{x}$ داده شده است تا با سریع تر به نتیجه برسد. همچنین دینامیک سیستم نیز در این قسمت تعریف شده است. تابع `render` نیز برای نمایش نتیجه به صورت شهودی تعریف شده است. کدهای این بخش به شکل زیر هستند:

```
class ThreeLinkRobotEnv(gym.Env):
    def __init__(self, render_mode = 'human'):
        super(ThreeLinkRobotEnv, self).__init__()

        # Define action and observation space
        self.action_space = spaces.Box(low=-0.1, high=0.1, shape=(3,),
dtype=np.float32)
        self.observation_space = spaces.Box(low=-np.pi, high=np.pi,
shape=(6,), dtype=np.float32)

        # Define robot parameters
        self.l1, self.l2, self.l3 = 1, 1, 1
        self.target = np.array([1, 0, 1, 0]) # Target position
        self.state = None

        # Reset the environment
        self.reset()

    def reset(self):
        # Set initial state close to the origin
        self.current_step = 0
        initial_theta = np.random.uniform(low=-np.pi/4, high=np.pi/4,
size=3)
        self.state = np.concatenate([initial_theta, np.zeros(3)])
        return self.state

    def step(self, action):

        self.current_step += 1
        theta1, theta2, theta3 = self.state[:3] + action
        x = (self.l1 * np.cos(theta1) + self.l2 * np.cos(theta1 +
theta2) + self.l3 * np.cos(theta1 + theta2 + theta3))
        y = (self.l1 * np.sin(theta1) + self.l2 * np.sin(theta1 +
theta2) + self.l3 * np.sin(theta1 + theta2 + theta3))

        distance_to_target = np.linalg.norm(self.target - np.array([x,
y]))
        reward = 20/distance_to_target # Negative distance to
encourage reaching the target
```

```

# Penalty for large actions to promote smooth movements
reward -= np.sum(np.abs(action)) * ۰,۳

self.state = np.array([theta۱, theta۲, theta۳, x, y,
distance_to_target])

done = distance_to_target < ۰,۱ or self.current_step >= ۲۰۰

return self.state, reward, done, {}

def render(self, render_mode="human"):
    theta۱, theta۲, theta۳ = self.state[:۳]
    x۱, y۱ = self.l۱ * np.cos(theta۱), self.l۱ * np.sin(theta۱)
    x۲, y۲ = x۱ + self.l۲ * np.cos(theta۱ + theta۲), y۱ + self.l۲ *
np.sin(theta۱ + theta۲)
    x۳, y۳ = x۲ + self.l۳ * np.cos(theta۱ + theta۲ + theta۳), y۲ +
self.l۳ * np.sin(theta۱ + theta۲ + theta۳)

    plt.figure()
    plt.plot([۰, x۱], [۰, y۱], 'ro-')
    plt.plot([x۱, x۲], [y۱, y۲], 'go-')
    plt.plot([x۲, x۳], [y۲, y۳], 'bo-')
    plt.plot(self.target[۰], self.target[۱], 'kx')
    plt.xlim(-۳, ۳)
    plt.ylim(-۳, ۳)
    plt.show()

```

۲/۳ یادگیری محیط

برای یادگیری محیط از الگوریتم یادگیری تقویتی PPO استفاده شده است. محیط با استفاده از این الگوریتم ۵ بار اجرا شده است. هر اپیزود نیز نمیتواند بیش از ۲۰۰ حرکت انجام دهد تا سرعت یادگیری بیشتر شود. برای این کار از کدهای زیر استفاده شده است:

```

env = TimeLimit(env, max_episode_steps=۲۰۰)
env = make_vec_env(lambda: env, n_envs=۱)

```

همچنین برای رسم نمودارها نیاز به مقادیر جایزه‌ها داریم. برای این کار، یک کلاس call back با نام RewardCallback تعریف شده است و یک شیء از آن به تابع learn داده شده است. به صورت زیر:

```

class RewardCallback(BaseCallback):
    def __init__(self, verbose=۰):
        super(RewardCallback, self).__init__(verbose)
        self.episode_rewards = []

    def _on_step(self) -> bool:
        if self.locals['dones'][۰]:

```

```

        self.episode_rewards.append(self.locals['rewards'][0])
    return True
reward_callback = RewardCallback()
callback = [reward_callback]
model.learn(total_timesteps=10000, callback=callback)
all_rewards.append(reward_callback.episode_rewards)

```

کدهای مربوط به بخش یادگیری نیز به صورت زیر هستند:

```

env = ThreeLinkRobotEnv()
env = TimeLimit(env, max_episode_steps=200)
env = make_vec_env(lambda: env, n_envs=1)

# Define the number of runs
num_runs = 5
all_rewards = []

# Run the training and collect rewards
for run in range(num_runs):
    model = PPO('MlpPolicy', env, verbose=1)

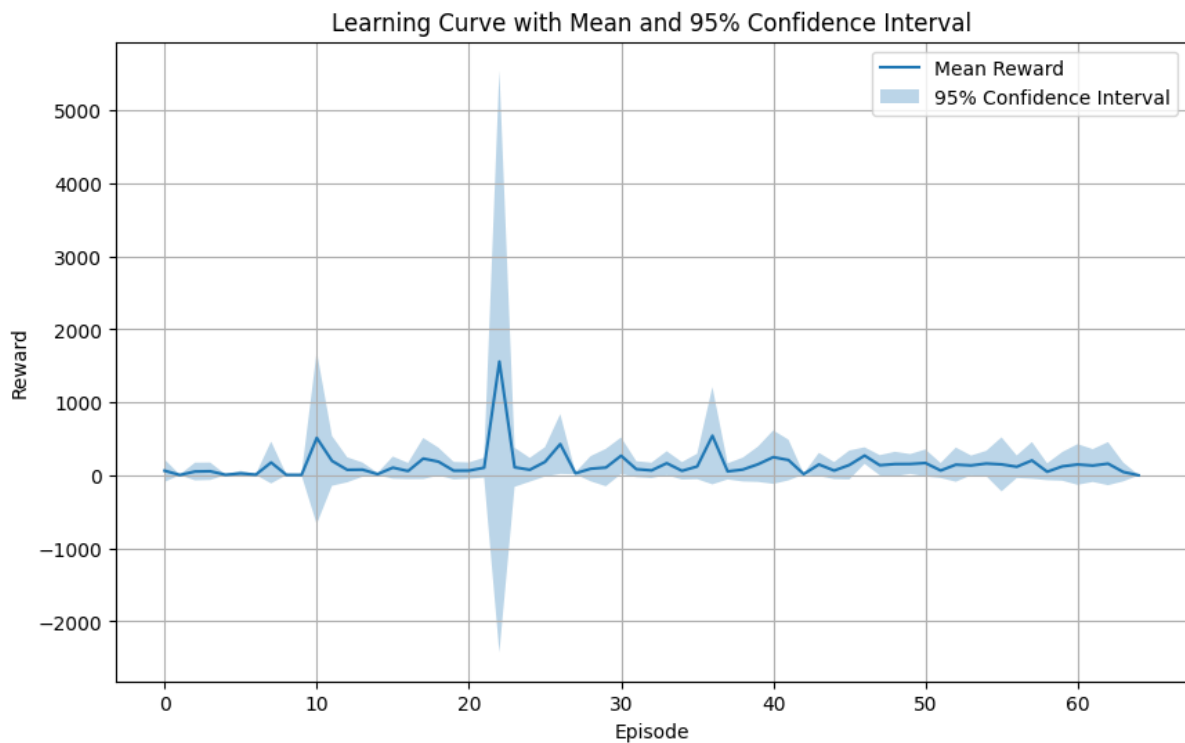
    reward_callback = RewardCallback()
    callback = [reward_callback]
    model.learn(total_timesteps=10000, callback=callback)
    all_rewards.append(reward_callback.episode_rewards)

# Print the collected rewards
for run in range(num_runs):
    print(f"Run {run+1} rewards: {all_rewards[run]}")

# Save the model
model.save("ppo_planar%dof")

```

نمودار Learning Curve with Mean and ۹۵% Confidence Interval برای یادگیری به صورت زیر به دست می‌آید:



۲/۴ رسم مدل

پس از یادگیری مدل، نمونه‌ی اجرا شده‌ی آن روی ربات مذکور را با کد زیر نمایش می‌دهیم:

```
model = PPO.load("ppo_planar3dof")

env = ThreeLinkRobotEnv()

# Test the trained model
obs = env.reset()
for i in range(1000):
    action, _states = model.predict(obs)
    obs, rewards, dones, _ = env.step(action)
    env.render()

# Stop the training if the target is reached
if dones:
    print("Target reached!")
    break
```

برخی از تصاویر به عنوان نمونه در زیر آورده خواهند شد:

