tags: `2024 年 下學期讀書計畫` `Reinforcement Learning`

# RL Homework 3: DDPG, TRPO, and PPO

## a

## ep_len



## ep_reward



## reward



## train

## hyperparameters

```python
num_episodes = 200
gamma = 0.995
tau = 0.002
hidden_size = 128
noise_scale = 0.3
replay_size = 100000
batch_size = 128
updates_per_step = 1
print_freq = 20
ewma_reward = 0
rewards = []
ewma_reward_history = []
total_numsteps = 0
updates = 0
```

## learning rates

```python
def __init__(self, num_inputs, action_space, gamma=0.995, tau=0.0005, hidden_size=128, lr_a=1e-4, lr_c=1e-3):
```

## NN architecture

```python
class Actor(nn.Module):
    def __init__(self, hidden_size, num_inputs, action_space):
        super(Actor, self).__init__()
        self.action_space = action_space
        num_outputs = action_space.shape[0]

        ########## YOUR CODE HERE (5~10 lines) ##########
        # Construct your own actor network
        self.fc1 = nn.Linear(num_inputs, hidden_size)
        self.fc2 = nn.Linear(hidden_size, hidden_size)
        self.fc3 = nn.Linear(hidden_size, num_outputs)
        self.relu = nn.ReLU()
        self.tanh = nn.Tanh()
```
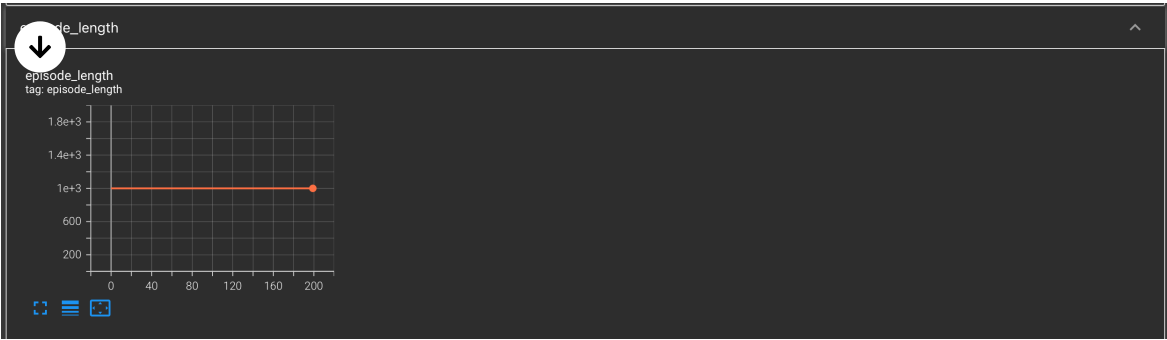
```python
class Critic(nn.Module):
    def __init__(self, hidden_size, num_inputs, action_space):
        super(Critic, self).__init__()
        self.action_space = action_space
        num_outputs = action_space.shape[0]

        ########## YOUR CODE HERE (5~10 lines) ##########
        # Construct your own critic network
        self.fc1 = nn.Linear(num_inputs + num_outputs, hidden_size)
        self.fc2 = nn.Linear(hidden_size, hidden_size)
        self.fc3 = nn.Linear(hidden_size, 1)
        self.relu = nn.ReLU()
```
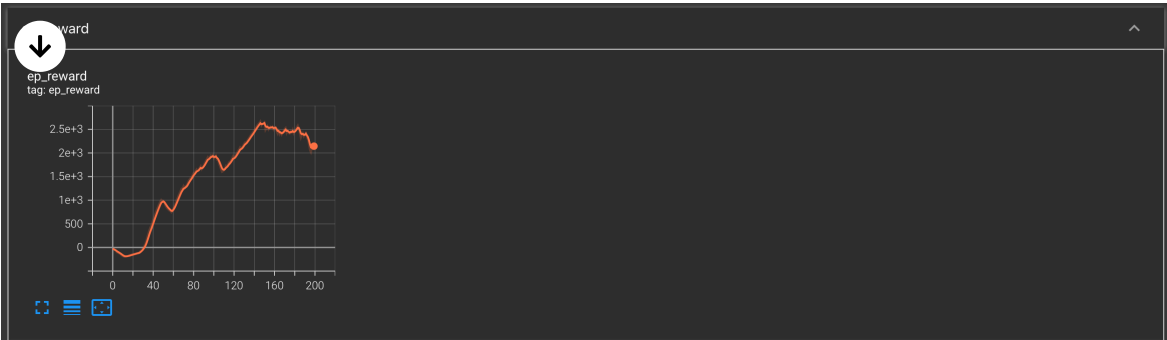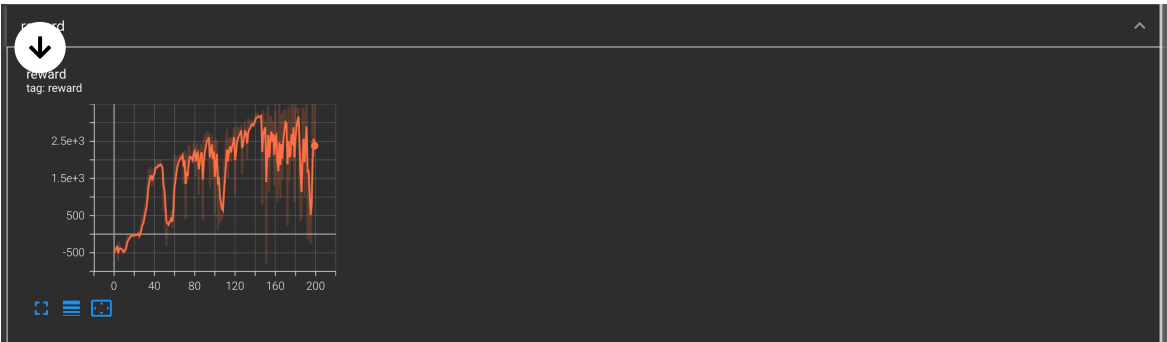
# b

## ep_len



## ep_reward



## reward

## train



## hyperparameters



```
num_episodes = 200
gamma = 0.995
tau = 0.002
hidden_size = 128
noise_scale = 0.3
replay_size = 100000
batch_size = 128
updates_per_step = 1
print_freq = 20
ewma_reward = 0
rewards = []
ewma_reward_history = []
total_numsteps = 0
updates = 0
```

## learning rates

```
def __init__(self, num_inputs, action_space, gamma=0.995, tau=0.0005, hidden_size=128, lr_a=1e-4, lr_c=1e-3):
```

## NN architecture

```python
class Actor(nn.Module):
    def __init__(self, hidden_size, num_inputs, action_space):
        super(Actor, self).__init__()
        self.action_space = action_space
        num_outputs = action_space.shape[0]

        ########## YOUR CODE HERE (5~10 lines) ##########
        # Construct your own actor network
        self.fc1 = nn.Linear(num_inputs, hidden_size)
        self.fc2 = nn.Linear(hidden_size, hidden_size)
        self.fc3 = nn.Linear(hidden_size, num_outputs)
        self.relu = nn.ReLU()
        self.tanh = nn.Tanh()
```

```python
class Critic(nn.Module):
    def __init__(self, hidden_size, num_inputs, action_space):
        super(Critic, self).__init__()
        self.action_space = action_space
        num_outputs = action_space.shape[0]

        ########## YOUR CODE HERE (5~10 lines) ##########
        # Construct your own critic network
        self.fc1 = nn.Linear(num_inputs + num_outputs, hidden_size)
        self.fc2 = nn.Linear(hidden_size, hidden_size)
        self.fc3 = nn.Linear(hidden_size, 1)
        self.relu = nn.ReLU()
```