## Implementation of Reinforce from scratch to play Cartpole-v1

- Environment documentations: https://gymnasium.farama.org/environments/classic_control/cart_pole/

## Install and import libraries

```
1 # Import Libraries
2
3 import gymnasium as gym
4 import torch
5 import torch.nn as nn
6 from torch.distributions.categorical import Categorical
7 from torch.distributions.normal import Normal
8 from collections import deque
9 import numpy as np
10
11 # For pushing to hub
12 from huggingface_hub import HfApi, snapshot_download
13 from huggingface_hub.repocard import metadata_eval_result, metadata_save
14 from huggingface_hub import notebook_login
15
16 from pathlib import Path
17 import datetime
18 import json
19 import imageio
20
21 import tempfile
22
23 import os
```

```
1 # Device
2
3 device = 'cuda' if torch.cuda.is_available() else 'cpu'
```

## 0. Visualise the observation space and action space

```
1 env = gym.make("CartPole-v1")
2
3 print(f"Number of available actions: {env.action_space.n}")
4 print(f"Sample a random action: {env.action_space.sample()}")
5
6 print(f"Sample a random observation: {env.observation_space.sample()}")
```

```
Number of available actions: 2
Sample a random action: 0
Sample a random observation: [4.745393   2.4669313  0.4147641  0.27528846]
```

## 1. Build a policy network using PyTorch

```
1 class Policy(nn.Module):
2   """
3   A policy network.
4
5   Args:
6     s_size (int): The size of 1 state space. \n
7     h_size (int): The number of hidden nodes in the network. \n
8     a_size (int): The number of distinct discrete actions, representing the number of output nodes. \n
9   """
10   def __init__(self, s_size, h_size, a_size):
11     super().__init__()
12
13     self.fc1 = nn.Sequential(nn.Linear(in_features = s_size,
14                                        out_features = h_size),
15                     nn.ReLU())
16
17     self.fc2 = nn.Sequential(nn.Linear(in_features = h_size,
18                                        out_features = a_size))
19
20   def forward(self, x):
21     """
22     The forward propagation of the policy network.
```

```
23
24     Args:
25       x (float tensor): Input to the network representing the observation / state, expected shape: (B, s_size). \n
26
27     Returns:
28       out(float tensor): Output of the network, representing the probability of taking each distinct discrete action, expecte shape:
29     """
30     out = self.fc1(x)
31     out = self.fc2(out)
32     out = torch.nn.functional.softmax(out, dim=-1)
33     return out
34
35   def act(self, state):
36     """
37     Sampling of an action.
38
39     Args:
40       state (float tensor): Input to the network representing the observation / state, expected shape: (B, s_size). \n
41
42     Returns:
43       action (int / int tensor): The index of the output nodes of the network, sampled based on the output probability of the networ
44       log_prob (float tensor): The ln of the probability of the action that was sampled based on the output probability of the newto
45
46     """
47     probs = self.forward(state).cpu()
48     m = Categorical(probs)
49     action = m.sample()
50     if action.shape[0] == 1:
51       return action.item(), m.log_prob(action)
52     else:
53       return action, m.log_prob(action)
54
55
```

## 2. Implementation of Reinforce algorithm

```
1  def reinforce(policy, optimizer, env, n_episodes, n_steps, device, gamma, print_every):
2
3    """
4    Train a policy network using the reinforce algorithm.
5
6    Args:
7      policy (nn.Module): A policy network. \n
8      optimizer (torch.optim): An optimizer. \n
9      env (gymnasium.env): An environment. \n
10     n_episodes (int): Number of training episodes. \n
11     n_steps (int): Maximum number of steps allowed in an episode. \n
12     device (str): 'cuda' or 'cpu'
13     gamma (float): A discount factor, range from 0 to 1.
14     print_every (int): Number of episode intervals to print the performance of the network. \n
15
16   Returns:
17     scores (list): A list of integer, each element representing the rewards scored in an episode. \n
18
19   """
20
21   # Create variables to store the rewards scored for every episode
22   scores = []
23
24   # This variable store up to rewards scored for every episode up to "print_every" episodes
25   scores_deque = deque(maxlen = print_every)
26
27
28   ####################
29   # For each episode #
30   ####################
31   for episode in range(1, n_episodes+1):
32
33     # Variable to store values for every step within an episode
34     reward_eps = [] # Reward scored by each step
35     log_prob_eps = [] # ln (prob of the action taken) for each step
36     returns_eps = deque(maxlen = n_steps) # discounted returns scored in each step
37     policy_loss_eps = [] # loss for each step -> ln (prob of the action taken) * discounted return
38
39     # Reset the environment for the beginning of each episode
40     state, info = env.reset()
41
42     ################
43     # For each step #
44     ################
```

```
45      for _ in range(n_steps):
46
47         # Sample an action using the policy network
48         action, log_prob = policy.act(torch.tensor(state).unsqueeze(0).to(device))
49
50         # Step the environment using the sampled action
51         state, reward, terminated, truncated, info = env.step(action)
52
53         # Store the rewards for this step and the ln (prob of this action)
54         reward_eps.append(reward)
55         log_prob_eps.append(log_prob)
56
57         # Check if this leads to termination or truncation
58         if terminated or truncated:
59            break
60
61      # Sum the rewards scored for the entire episode and store them
62      scores.append(sum(reward_eps))
63      scores_deque.append(sum(reward_eps))
64
65      # Calculate the discounted returns
66      for t in range(len(reward_eps))[::-1]:
67         returns = reward_eps[t] + gamma * returns_eps[0] if len(returns_eps) > 0 else reward_eps[t]
68         returns_eps.appendleft(returns)
69
70      # Normalize the discounted returns
71      eps = np.finfo(np.float32).eps.item() # eps is smallest reprsentatable float
72      returns_eps = torch.tensor(returns_eps) # Convert into torch tensor for later calculation
73      returns_eps = (returns_eps - returns_eps.mean()) / (returns_eps.std() + eps)
74
75      # Calculate the loss
76      for log_prob, returns in zip(log_prob_eps, returns_eps):
77         policy_loss_eps.append(-log_prob * returns) # torch tensor * torch tensor
78      loss = torch.cat(policy_loss_eps).sum() # cat -> makes into one torch tensor, sum() -> summation
79
80      # Backward propagation and gradient descent
81      optimizer.zero_grad()
82      loss.backward()
83      optimizer.step()
84
85      # Print information
86      if episode % print_every == 0:
87         print(f"Current Episode: {episode} | Average reward: {np.mean(scores_deque)}")
88
89   return scores
90
91
92
93
94
95
```

## 3. Train the policy network using reinforce algorithm

```
1 #0. Create device
2 device = 'cuda' if torch.cuda.is_available() else 'cpu'
3
4 # 1. Create hyperparameters
5 cartpole_hyperparameters = {
6     "h_size": 16,
7     "n_training_episodes": 1000,
8     "n_evaluation_episodes": 10,
9     "max_t": 1000,
10    "gamma": 1.0,
11    "lr": 0.01,
12    "env_id": 'CartPole-v1',
13    "state_space": 4,
14    "action_space": 2,
15 }
16
17 # 2. Create environment
18 env = gym.make(cartpole_hyperparameters['env_id'])
19
20 # 3. Create policy network
21 cartpole_policy = Policy(cartpole_hyperparameters['state_space'],
22                          cartpole_hyperparameters['h_size'],
23                          cartpole_hyperparameters['action_space']).to(device)
24
25 # 4. Create optimizer
26 optimizer = torch.optim.Adam(cartpole_policy.parameters(),
```

```
27                          lr = cartpole_hyperparameters['lr'])
28
29 # 5. Training loop
30 scores = reinforce(policy = cartpole_policy,
31                    optimizer = optimizer,
32                    env = env,
33                    n_episodes = cartpole_hyperparameters['n_training_episodes'],
34                    n_steps = cartpole_hyperparameters['max_t'],
35                    device = device,
36                    gamma = cartpole_hyperparameters['gamma'],
37                    print_every = 100)
38
```

```
Current Episode: 100 | Average reward: 44.33
Current Episode: 200 | Average reward: 308.61
Current Episode: 300 | Average reward: 322.93
Current Episode: 400 | Average reward: 388.92
Current Episode: 500 | Average reward: 352.5
Current Episode: 600 | Average reward: 459.3
Current Episode: 700 | Average reward: 476.28
Current Episode: 800 | Average reward: 500.0
Current Episode: 900 | Average reward: 461.11
Current Episode: 1000 | Average reward: 379.94
```

## ∨ 4. Evaluate the agent

```
1 def evaluate_agent(n_eval_episodes, n_steps, policy, env, device):
2
3   """
4   Evaluate the performance of an agent by calculating the mean and standard deviation rewards over n_eval_episodes of episodes.
5
6   Args:
7     n_eval_episodes (int): Number of evaluation episodes. \n
8     n_steps (int): Maximum number of steps allowed in an episode. \n
9     policy (nn.Module): A policy network. \n
10    env (gymnasium.env): An environment. \n
11    device (str): 'cuda' or 'cpu'. \n
12
13  Returns:
14    mean_reward (float): Mean reward scored across the evaluated episodes.
15    std_reward (float): Standard deviation reward scored across the evaluated episodes.
16  """
17
18
19  rewards_across_episodes = [] # Contains rewards scored in each episode
20
21  ####################
22  # For each episode #
23  ####################
24  for episode in range(n_eval_episodes):
25
26    # To store reward scored in each step
27    rewards = []
28
29    # Reset the environment
30    state, info = env.reset()
31
32    #################
33    # For each step #
34    #################
35    for step in range(n_steps):
36
37      # Sample an action
38      action, _ = policy.act(torch.tensor(state).unsqueeze(0).to(device))
39
40      # Step the environment by taking the action
41      state, reward, terminated, truncated, info = env.step(action)
42
43      # Store the reward scored in this step
44      rewards.append(reward)
45
46      # Check if truncated or terminated
47      if truncated or terminated:
48        break
49
50    # Sum the reward scored in the entire episode and store it
51    rewards_across_episodes.append(sum(rewards))
52
53  # Calculate the mean and standard deviation
54  mean_reward = np.array(rewards_across_episodes).mean()
55  std_reward = np.array(rewards_across_episodes).std()
```

```
56    return mean_reward, std_reward
57
```

```
1 mean_reward, std_reward = evaluate_agent(n_eval_episodes = cartpole_hyperparameters['n_evaluation_episodes'],
2                                          n_steps = cartpole_hyperparameters['max_t'],
3                                          policy = cartpole_policy,
4                                          env = env,
5                                          device = device)
6
7 print(f"Mean reward: {mean_reward}, standard deviation: {std_reward}")
```

⤷ Mean reward: 500.0, standard deviation: 0.0

## ⌄ 5. Push to Hub

- Code source: https://colab.research.google.com/github/wengti/Reinforcement-Learning-Tutorial-/blob/main/notebooks/unit4/unit4.ipynb#scrollTo=LIVsvlW_8tcw

- Creat a write token here: https://huggingface.co/settings/tokens/new?tokenType=write

```
1 def record_video(env, policy, out_directory, device, fps=30):
2   """
3   Generate a replay video of the agent
4   :param env
5   :param Qtable: Qtable of our agent
6   :param out_directory
7   :param fps: how many frame per seconds (with taxi-v3 and frozenlake-v1 we use 1)
8   """
9   images = []
10  state, info = env.reset()
11  terminated = False
12  truncated = False
13  img = env.render()
14  images.append(img)
15  while not terminated and not truncated:
16    # Take the action (index) that have the maximum expected future reward given that state
17    action, _ = policy.act(torch.tensor(state).unsqueeze(0).to(device))
18    state, reward, terminated, truncated, info = env.step(action) # We directly put next_state = state for recording logic
19    img = env.render()
20    images.append(img)
21  imageio.mimsave(out_directory, [np.array(img) for i, img in enumerate(images)], fps=fps)
```

```
1 def push_to_hub(repo_id,
2                 model,
3                 hyperparameters,
4                 eval_env,
5                 video_fps=30
6                 ):
7   """
8   Evaluate, Generate a video and Upload a model to Hugging Face Hub.
9   This method does the complete pipeline:
10  - It evaluates the model
11  - It generates the model card
12  - It generates a replay video of the agent
13  - It pushes everything to the Hub
14
15  :param repo_id: repo_id: id of the model repository from the Hugging Face Hub
16  :param model: the pytorch model we want to save
17  :param hyperparameters: training hyperparameters
18  :param eval_env: evaluation environment
19  :param video_fps: how many frame per seconds to record our video replay
20  """
21
22  _, repo_name = repo_id.split("/")
23  api = HfApi()
24
25  # Step 1: Create the repo
26  repo_url = api.create_repo(
27        repo_id=repo_id,
28        exist_ok=True,
29  )
30
31  with tempfile.TemporaryDirectory() as tmpdirname:
32    local_directory = Path(tmpdirname)
33
34    # Step 2: Save the model
35    torch.save(model, local_directory / "model.pt")
36
```

```python
37     # Step 3: Save the hyperparameters to JSON
38     with open(local_directory / "hyperparameters.json", "w") as outfile:
39       json.dump(hyperparameters, outfile)
40
41     # Step 4: Evaluate the model and build JSON
42     mean_reward, std_reward = evaluate_agent(hyperparameters["n_evaluation_episodes"],
43                                              hyperparameters["max_t"],
44                                              model,
45                                              eval_env,
46                                              'cuda')
47     # Get datetime
48     eval_datetime = datetime.datetime.now()
49     eval_form_datetime = eval_datetime.isoformat()
50
51     evaluate_data = {
52         "env_id": hyperparameters["env_id"],
53         "mean_reward": mean_reward,
54         "n_evaluation_episodes": hyperparameters["n_evaluation_episodes"],
55         "eval_datetime": eval_form_datetime,
56     }
57
58     # Write a JSON file
59     with open(local_directory / "results.json", "w") as outfile:
60         json.dump(evaluate_data, outfile)
61
62     # Step 5: Create the model card
63     env_name = hyperparameters["env_id"]
64     env_id = env_name
65
66     metadata = {}
67     metadata["tags"] = [
68         env_name,
69         "reinforce",
70         "reinforcement-learning",
71         "custom-implementation",
72         "deep-rl-class"
73     ]
74
75     # Add metrics
76     eval = metadata_eval_result(
77         model_pretty_name=repo_name,
78         task_pretty_name="reinforcement-learning",
79         task_id="reinforcement-learning",
80         metrics_pretty_name="mean_reward",
81         metrics_id="mean_reward",
82         metrics_value=f"{mean_reward:.2f} +/- {std_reward:.2f}",
83         dataset_pretty_name=env_name,
84         dataset_id=env_name,
85     )
86
87     # Merges both dictionaries
88     metadata = {**metadata, **eval}
89
90     model_card = f"""
91 # **Reinforce** Agent playing **{env_id}**
92 This is a trained model of a **Reinforce** agent playing **{env_id}** .
93 To learn to use this model and train yours check Unit 4 of the Deep Reinforcement Learning Course: https://huggingface.co/deep-rl
94     """
95
96     readme_path = local_directory / "README.md"
97     readme = ""
98     if readme_path.exists():
99         with readme_path.open("r", encoding="utf8") as f:
100            readme = f.read()
101    else:
102      readme = model_card
103
104    with readme_path.open("w", encoding="utf-8") as f:
105      f.write(readme)
106
107    # Save our metrics to Readme metadata
108    metadata_save(readme_path, metadata)
109
110    # Step 6: Record a video
111    video_path =  local_directory / "replay.mp4"
112    record_video(eval_env, model, video_path, 'cuda', video_fps)
113
114    # Step 7. Push everything to the Hub
115    api.upload_folder(
116        repo_id=repo_id,
117        folder_path=local_directory,
118        path_in_repo=".",
```

```
119       )
120
121       print(f"Your model is pushed to the Hub. You can view your model here: {repo_url}")
```

```
1 # Login to hugging face with a write token
2
3 notebook_login()
```

```
1 # Create a new environment with render mode (needed for video recording)
2 eval_env = gym.make(cartpole_hyperparameters['env_id'], render_mode = 'rgb_array')
3
4 # Push to Hub
5 push_to_hub(repo_id = "wengti0608/Reinforce-Cartpole-v1-attempt2",
6             model = cartpole_policy,
7             hyperparameters = cartpole_hyperparameters,
8             eval_env = eval_env)
```

```
WARNING:imageio_ffmpeg:IMAGEIO FFMPEG_WRITER WARNING: input image is not divisible by macro_block_size=16, resizing from (600, 400)
Uploading...: 100%                                          3.78k/3.78k [00:01<00:00, 18.8kB/s]

Your model is pushed to the Hub. You can view your model here: https://huggingface.co/wengti0608/Reinforce-Cartpole-v1-attempt2
```

## Practice: Application of Reinforce to play Continuous Mountain Car

- Environment documentation: https://gymnasium.farama.org/environments/classic_control/mountain_car/
- Continuous Mountain Car agent takes continuous actions. Therefore, the policy network is slightly modified to output mean and standard deviation instead of index of discrete actions.

## 0. Visualize Environment

```
1 # Visualize the environment
2
3 env = gym.make("MountainCarContinuous-v0")
4
5 print(f"Randomly sample an action: {env.action_space.sample()}")
6 print(f"Randomly sample a state: {env.observation_space.sample()}")
7
```

```
Randomly sample an action: [0.9539736]
Randomly sample a state: [-0.00267963 -0.00615319]
```

## 1. Create the policy network

```
1 class Policy(nn.Module):
2   """
3   A policy network.
4
5   Args:
6     s_size (int): The size of 1 state space. \n
7     h_size (int): The number of hidden nodes in the network. \n
8     a_size (int): Number of continuos actions. \n
9   """
10   def __init__(self, s_size, h_size, a_size):
11     super().__init__()
12
13     self.fc1 = nn.Sequential(nn.Linear(in_features = s_size,
14                                        out_features = h_size),
15                              nn.ReLU())
16
17     self.fc2 = nn.Sequential(nn.Linear(in_features = h_size,
18                                        out_features = h_size),
19                              nn.ReLU())
20
21     self.mean_head = nn.Linear(in_features = h_size,
22                                out_features = a_size)
23
24     self.log_std_head = nn.Linear(in_features = h_size,
25                                   out_features = a_size)
26
27   def forward(self, x):
```

```
28     """
29     The forward propagation of the policy network.
30
31     Args:
32       x (float tensor): Input to the network representing the observation / state, expected shape: (B, s_size). \n
33
34     Returns:
35       mean (float tensor): Mean of a Normal Distribution, (B, a_size). \n
36       std (float tensor): Standard deviation of a Normal Distribtuion, (B, a_size). \n
37     """
38     out = self.fc1(x)
39     out = self.fc2(out)
40
41     mean = self.mean_head(out)
42     log_std = self.log_std_head(out)
43     std = torch.exp(log_std)
44
45     return mean, std
46
47   def act(self, state):
48     """
49     Sampling of an action.
50
51     Args:
52       state (float tensor): Input to the network representing the observation / state, expected shape: (B, s_size). \n
53
54     Returns:
55       action_clipped (float): The value of the action taken, in the range of -1 to 1 \n
56       log_prob (float tensor): The ln of the probability of the action that was sampled based on the output probability of the newto
57
58     """
59     mean, std = self.forward(state)
60     mean = mean.cpu() # torch.tensor, (1,1)
61     std = std.cpu() # torch.tensor, (1,1)
62
63     m = Normal(mean, std)
64
65     action = m.sample() # torch.tensor, (1,1)
66     action_clipped = torch.clamp(action, -1, 1).item() #float
67
68     log_prob = m.log_prob(action)[0] # torch.tensor, (1,)
69
70     return action_clipped, log_prob
71
72
```

## 2. Implement reinforce algorithm from scratch

```
1  def reinforce(policy, optimizer, env, n_episodes, n_steps, device, gamma, print_every):
2
3    """
4    Train a policy network using the reinforce algorithm.
5
6    Args:
7      policy (nn.Module): A policy network. \n
8      optimizer (torch.optim): An optimizer. \n
9      env (gymnasium.env): An environment. \n
10     n_episodes (int): Number of training episodes. \n
11     n_steps (int): Maximum number of steps allowed in an episode. \n
12     device (str): 'cuda' or 'cpu'
13     gamma (float): A discount factor, range from 0 to 1.
14     print_every (int): Number of episode intervals to print the performance of the network. \n
15
16   Returns:
17     scores (list): A list of integer, each element representing the rewards scored in an episode. \n
18
19   """
20
21   # Create variables to store the rewards scored for every episode
22   scores = []
23
24   # This variable store up to rewards scored for every episode up to "print_every" episodes
25   scores_deque = deque(maxlen = print_every)
26   len_deque = deque(maxlen = print_every)
27
28
29   ####################
30   # For each episode #
31   ####################
32   for episode in range(1, n_episodes+1):
```

```
33
34      # Variable to store values for every step within an episode
35      reward_eps = [] # Reward scored by each step
36      log_prob_eps = [] # ln (prob of the action taken) for each step
37      returns_eps = deque(maxlen = n_steps) # discounted returns scored in each step
38      policy_loss_eps = [] # loss for each step -> ln (prob of the action taken) * discounted return
39
40      # Reset the environment for the beginning of each episode
41      state, info = env.reset()
42
43      #################
44      # For each step #
45      #################
46      for _ in range(n_steps):
47
48        # Sample an action using the policy network
49        action, log_prob = policy.act(torch.tensor(state).unsqueeze(0).to(device))
50
51        # Step the environment using the sampled action
52        state, reward, terminated, truncated, info = env.step(np.array([action]))
53
54        # Store the rewards for this step and the ln (prob of this action)
55        reward_eps.append(reward)
56        log_prob_eps.append(log_prob)
57
58        # Check if this leads to termination or truncation
59        if terminated or truncated:
60          break
61
62      # Sum the rewards scored for the entire episode and store them
63      scores.append(sum(reward_eps))
64      scores_deque.append(sum(reward_eps))
65      len_deque.append(len(reward_eps))
66
67      # Calculate the discounted returns
68      for t in range(len(reward_eps))[::-1]:
69        returns = reward_eps[t] + gamma * returns_eps[0] if len(returns_eps) > 0 else reward_eps[t]
70        returns_eps.appendleft(returns)
71
72      # Normalize the discounted returns
73      eps = np.finfo(np.float32).eps.item() # eps is smallest reprsentatable float
74      returns_eps = torch.tensor(returns_eps) # Convert into torch tensor for later calculation
75      returns_eps = (returns_eps - returns_eps.mean()) / (returns_eps.std() + eps)
76
77      # Calculate the loss
78      for log_prob, returns in zip(log_prob_eps, returns_eps):
79        policy_loss_eps.append(-log_prob * returns) # torch tensor * torch tensor
80      loss = torch.cat(policy_loss_eps).sum() # cat -> makes into one torch tensor, sum() -> summation
81
82      # Backward propagation and gradient descent
83      optimizer.zero_grad()
84      loss.backward()
85      optimizer.step()
86
87      # Print information
88      if episode % print_every == 0:
89        print(f"Current Episode: {episode} | Average reward: {np.mean(scores_deque)} | Average length of ep: {np.mean(len_deque)}")
90
91    return scores
92
93
94
95
96
97
```

## 3. Training

```
1 # 0. Device
2 device = 'cuda' if torch.cuda.is_available() else 'cpu'
3
4 # 1. Hyperparameters
5 car_hyperparameters = {
6     "h_size": 64,
7     "n_training_episodes": 2000,
8     "n_evaluation_episodes": 10,
9     "max_t": 999,
10    "gamma": 0.9,
11    "lr": 0.00001,
12    "env_id": 'MountainCarContinuous-v0',
```

```
13      "state_space": 2,
14      "action_space": 1,
15 }
16
17 # 2. Build a Policy Network
18 car_policy = Policy(car_hyperparameters['state_space'],
19                     car_hyperparameters['h_size'],
20                     car_hyperparameters['action_space']).to(device)
21
22 # 3. Create environment
23 env = gym.make(car_hyperparameters['env_id'])
24
25 # 4. Train the network
26 optimizer = torch.optim.Adam(car_policy.parameters(),
27                              lr = car_hyperparameters['lr'])
28
29 scores = reinforce(policy = car_policy,
30                    optimizer = optimizer,
31                    env = env,
32                    n_episodes = car_hyperparameters['n_training_episodes'],
33                    n_steps = car_hyperparameters['max_t'],
34                    device = device,
35                    gamma = car_hyperparameters['gamma'],
36                    print_every = 100)
37
```

⮂  Show hidden output

## ⌄ 4. Evaluation

```
 1 def evaluate_agent(n_eval_episodes, n_steps, policy, env, device):
 2
 3   """
 4   Evaluate the performance of an agent by calculating the mean and standard deviation rewards over n_eval_episodes of episodes.
 5
 6   Args:
 7     n_eval_episodes (int): Number of evaluation episodes. \n
 8     n_steps (int): Maximum number of steps allowed in an episode. \n
 9     policy (nn.Module): A policy network. \n
10     env (gymnasium.env): An environment. \n
11     device (str): 'cuda' or 'cpu'. \n
12
13   Returns:
14     mean_reward (float): Mean reward scored across the evaluated episodes.
15     std_reward (float): Standard deviation reward scored across the evaluated episodes.
16   """
17
18
19   rewards_across_episodes = [] # Contains rewards scored in each episode
20
21   ####################
22   # For each episode #
23   ####################
24   for episode in range(n_eval_episodes):
25
26     # To store reward scored in each step
27     rewards = []
28
29     # Reset the environment
30     state, info = env.reset()
31
32     #################
33     # For each step #
34     #################
35     for step in range(n_steps):
36
37       # Sample an action
38       action, _ = policy.act(torch.tensor(state).unsqueeze(0).to(device))
39
40       # Step the environment by taking the action
41       state, reward, terminated, truncated, info = env.step(np.array([action]))
42
43       # Store the reward scored in this step
44       rewards.append(reward)
45
46       # Check if truncated or terminated
47       if truncated or terminated:
48         break
49
50     # Sum the reward scored in the entire episode and store it
```

```
51    rewards_across_episodes.append(sum(rewards))
52
53  # Calculate the mean and standard deviation
54  mean_reward = np.array(rewards_across_episodes).mean()
55  std_reward = np.array(rewards_across_episodes).std()
56  return mean_reward, std_reward
57
```

```
1 mean_reward, std_reward = evaluate_agent(n_eval_episodes = car_hyperparameters['n_evaluation_episodes'],
2                                         n_steps = car_hyperparameters['max_t'],
3                                         policy = car_policy,
4                                         env = env,
5                                         device = device)
6
7 print(f"Mean reward: {mean_reward}, Std reward: {std_reward}")
```

```
→  Mean reward: -46.63831101175287, Std reward: 1.2278728245384134
```

## ⌄ 5. Push to Hub

```
1 def record_video(env, policy, out_directory, device, fps=30):
2   """
3   Generate a replay video of the agent
4   :param env
5   :param Qtable: Qtable of our agent
6   :param out_directory
7   :param fps: how many frame per seconds (with taxi-v3 and frozenlake-v1 we use 1)
8   """
9   images = []
10   state, info = env.reset()
11   terminated = False
12   truncated = False
13   img = env.render()
14   images.append(img)
15   while not terminated and not truncated:
16     # Take the action (index) that have the maximum expected future reward given that state
17     action, _ = policy.act(torch.tensor(state).unsqueeze(0).to(device))
18     state, reward, terminated, truncated, info = env.step(np.array([action])) # We directly put next_state = state for recording log
19     img = env.render()
20     images.append(img)
21   imageio.mimsave(out_directory, [np.array(img) for i, img in enumerate(images)], fps=fps)
```

```
1 def push_to_hub(repo_id,
2                 model,
3                 hyperparameters,
4                 eval_env,
5                 video_fps=30
6                 ):
7   """
8   Evaluate, Generate a video and Upload a model to Hugging Face Hub.
9   This method does the complete pipeline:
10   - It evaluates the model
11   - It generates the model card
12   - It generates a replay video of the agent
13   - It pushes everything to the Hub
14
15   :param repo_id: repo_id: id of the model repository from the Hugging Face Hub
16   :param model: the pytorch model we want to save
17   :param hyperparameters: training hyperparameters
18   :param eval_env: evaluation environment
19   :param video_fps: how many frame per seconds to record our video replay
20   """
21
22   _, repo_name = repo_id.split("/")
23   api = HfApi()
24
25   # Step 1: Create the repo
26   repo_url = api.create_repo(
27       repo_id=repo_id,
28       exist_ok=True,
29   )
30
31   with tempfile.TemporaryDirectory() as tmpdirname:
32     local_directory = Path(tmpdirname)
33
34     # Step 2: Save the model
35     torch.save(model, local_directory / "model.pt")
36
```

```python
37      # Step 3: Save the hyperparameters to JSON
38      with open(local_directory / "hyperparameters.json", "w") as outfile:
39          json.dump(hyperparameters, outfile)
40
41      # Step 4: Evaluate the model and build JSON
42      mean_reward, std_reward = evaluate_agent(hyperparameters["n_evaluation_episodes"],
43                                                hyperparameters["max_t"],
44                                                model,
45                                                eval_env,
46                                                'cuda')
47      # Get datetime
48      eval_datetime = datetime.datetime.now()
49      eval_form_datetime = eval_datetime.isoformat()
50
51      evaluate_data = {
52          "env_id": hyperparameters["env_id"],
53          "mean_reward": mean_reward,
54          "n_evaluation_episodes": hyperparameters["n_evaluation_episodes"],
55          "eval_datetime": eval_form_datetime,
56      }
57
58      # Write a JSON file
59      with open(local_directory / "results.json", "w") as outfile:
60          json.dump(evaluate_data, outfile)
61
62      # Step 5: Create the model card
63      env_name = hyperparameters["env_id"]
64      env_id = env_name
65
66      metadata = {}
67      metadata["tags"] = [
68          env_name,
69          "reinforce",
70          "reinforcement-learning",
71          "custom-implementation",
72          "deep-rl-class"
73      ]
74
75      # Add metrics
76      eval = metadata_eval_result(
77          model_pretty_name=repo_name,
78          task_pretty_name="reinforcement-learning",
79          task_id="reinforcement-learning",
80          metrics_pretty_name="mean_reward",
81          metrics_id="mean_reward",
82          metrics_value=f"{mean_reward:.2f} +/- {std_reward:.2f}",
83          dataset_pretty_name=env_name,
84          dataset_id=env_name,
85      )
86
87      # Merges both dictionaries
88      metadata = {**metadata, **eval}
89
90      model_card = f"""
91 # **Reinforce** Agent playing **{env_id}**
92 This is a trained model of a **Reinforce** agent playing **{env_id}** .
93 To learn to use this model and train yours check Unit 4 of the Deep Reinforcement Learning Course: https://huggingface.co/deep-r]
94      """
95
96      readme_path = local_directory / "README.md"
97      readme = ""
98      if readme_path.exists():
99          with readme_path.open("r", encoding="utf8") as f:
100             readme = f.read()
101     else:
102         readme = model_card
103
104     with readme_path.open("w", encoding="utf-8") as f:
105         f.write(readme)
106
107     # Save our metrics to Readme metadata
108     metadata_save(readme_path, metadata)
109
110     # Step 6: Record a video
111     video_path =  local_directory / "replay.mp4"
112     record_video(eval_env, model, video_path, 'cuda', video_fps)
113
114     # Step 7. Push everything to the Hub
115     api.upload_folder(
116         repo_id=repo_id,
117         folder_path=local_directory,
118         path_in_repo=".",
```

```
119    )
120
121    print(f"Your model is pushed to the Hub. You can view your model here: {repo_url}")
```

```
1 # Login to HuggingFace Hub
2 notebook_login()
```

```
1 # Create an evaluation environment
2 eval_env = gym.make(car_hyperparameters['env_id'], render_mode = 'rgb_array')
3
4 # Push to Hub
5 push_to_hub(repo_id = "wengti0608/Reinforce-MountainCarContinuous-v0-attempt1",
6            model = car_policy,
7            hyperparameters = car_hyperparameters,
8            eval_env = eval_env)
```

> WARNING:imageio_ffmpeg:IMAGEIO FFMPEG_WRITER WARNING: input image is not divisible by macro_block_size=16, resizing from (600, 400)
>
> Uploading...: 100%                                              127k/127k [00:01<00:00, 637kB/s]
>
> Your model is pushed to the Hub. You can view your model here: https://huggingface.co/wengti0608/Reinforce-MountainCarContinuous-v0

## ⌄ Solving Mountain Car with SAC

- Reinforce did not manage to solve Continuous Mountain Car. Therefore, SAC is used instead.

```
1 !pip install stable-baselines3==2.0.0a5
```

> Show hidden output

```
1 !pip install huggingface_sb3
```

> Show hidden output

## ⌄ 1. Training

```
1 import gymnasium as gym
2 from stable_baselines3.sac import SAC
3 from stable_baselines3.common.callbacks import EvalCallback
4 from stable_baselines3.common.monitor import Monitor
5 from stable_baselines3.common.evaluation import evaluate_policy
6
7
8 # 1. Make environment
9 env = gym.make("MountainCarContinuous-v0")
10
11 # 2. Create a SAC model
12
13 # The hyperparameter is provided here:
14 # https://huggingface.co/sb3/sac-MountainCarContinuous-v0
15 policy_kwargs = {'log_std_init': -3.67,
16                  'net_arch': [64, 64]}
17
18 model = SAC(batch_size = 512,
19            buffer_size = 50000,
20            ent_coef = 0.1,
21            gamma = 0.9999,
22            gradient_steps = 32,
23            learning_rate = 0.0003,
24            learning_starts = 0,
25            policy = 'MlpPolicy',
26            policy_kwargs = policy_kwargs,
27            tau = 0.01,
28            train_freq = 32,
29            use_sde = True,
30            env = env)
31
32 # 3. Train the model
33
34 # As SAC does not output rollout on its own
35 # A callback is manually created...
36 eval_env = Monitor(gym.make("MountainCarContinuous-v0", render_mode = "rgb_array"))
37
38 eval_callback = EvalCallback(
39     eval_env,
```

```
40      best_model_save_path = './logs/SAC',
41      log_path = './logs/SAC',
42      eval_freq = 1e3,
43      deterministic = True,
44      render = False
45 )
46
47 # Training
48 model.learn(total_timesteps = 5e4, callback = eval_callback)
49
50 # 4. Save the model
51 model_name = "SAC-MountainCarContinuous-v0"
52 model.save(model_name)
53
```

Show hidden output

## 2. Evaluation

```
 1 #1. Create an evaluation environment
 2 eval_env = Monitor(gym.make("MountainCarContinuous-v0", render_mode = "rgb_array"))
 3
 4 #2. Evaluate policy
 5 mean_reward, std_reward = evaluate_policy(model = model,
 6                                           env = eval_env,
 7                                           n_eval_episodes = 10,
 8                                           deterministic = True)
 9
10 print(f"The mean reward: {mean_reward} | The standard deviation: {std_reward}")
```

The mean reward: 94.67033819999999 | The standard deviation: 0.2572431881083725

## 3. Push to Hub

```
1 from huggingface_hub import notebook_login
2
3 notebook_login()
```

```
 1 from stable_baselines3.common.vec_env import DummyVecEnv
 2 from huggingface_sb3 import package_to_hub
 3
 4 package_to_hub(model = model,
 5                model_name = model_name,
 6                model_architecture = "SAC",
 7                env_id = "MountainCarContinuous-v0",
 8                eval_env = DummyVecEnv([lambda: Monitor(gym.make("MountainCarContinuous-v0", render_mode = "rgb_array"))]),
 9                repo_id = "wengti0608/SAC-MountainCarContinuous-v0",
10                commit_message = "First Commit")
```

Show hidden output