

Assignment 2: Deep Q Learning and Policy Gradient

CS260R 2023Fall: Reinforcement Learning. Department of Computer Science at University of California, Los Angeles. Course Instructor: Professor Bolei ZHOU. Assignment author: Zhenghao PENG, Yiran WANG.

Student Name	Student ID
Shaira Alam	506302126

Welcome to the assignment 2 of our RL course. This assignment consists of three parts:

- Section 2: Implement Q learning in tabular setting (20 points)
- Section 3: Implement Deep Q Network with pytorch (30 points)
- Section 4: Implement policy gradient method REINFORCE with pytorch (30 points)
- Section 5: Implement policy gradient method with baseline (20 points) (+20 points bonus)

Section 0 and Section 1 set up the dependencies and prepare some useful functions.

The experiments we'll conduct and their expected goals:

1. Naive Q learning in FrozenLake (should solve)
2. DQN in CartPole (should solve)
3. DQN in MetaDrive-Easy (should solve)
4. Policy Gradient w/o baseline in CartPole (w/ and w/o advantage normalization) (should solve)
5. Policy Gradient w/o baseline in MetaDrive-Easy (should solve)
6. Policy Gradient w/ baseline in CartPole (w/ advantage normalization) (should solve)
7. Policy Gradient w/ baseline in MetaDrive-Easy (should solve)
8. Policy Gradient w/ baseline in MetaDrive-Hard (>20 return) (Optional, +20 points bonus can be earned)

NOTE: MetaDrive does not support python=3.12. If you are in python=3.12, we suggest to recreate a new conda environment:

```
conda env remove -n cs260r
conda create -n cs260r python=3.11 -y
pip install notebook # Install jupyter notebook
jupyter notebook # Run jupyter notebook
```

Section 0: Dependencies

Please install the following dependencies.

Notes on MetaDrive

MetaDrive is a lightweight driving simulator which we will use for DQN and Policy Gradient methods. It can not be run on M1-chip Mac. We suggest using Colab or Linux for running MetaDrive.

Please ignore this warning from MetaDrive: `WARNING:root:BaseEngine is not launched, fail to sync seed to engine!`

Notes on Colab

We have several cells used for installing dependencies for Colab only. Please make sure they are run properly.

You don't need to install python packages again and again after **restarting the runtime**, since the Colab instance still remembers the python environment after you installing packages for the first time. But you do need to rerun those packages installation script after you **reconnecting to the runtime** (which means Google assigns a new machine to you and thus the python environment is new).

```
In [1]: RUNNING_IN_COLAB = 'google.colab' in str(get_ipython()) # Detect if it is running in Colab
```

```
In [2]: # Similar to AS1

!pip install -U pip
!pip install numpy scipy "gymnasium<0.29"
!pip install torch torchvision
!pip install mediapy
```

```
Requirement already satisfied: pip in c:\users\user\anaconda3\lib\site-packages (22.2.2)
Collecting pip
  Using cached pip-23.3.1-py3-none-any.whl (2.1 MB)
ERROR: To modify pip, please run the following command:
C:\Users\User\anaconda3\python.exe -m pip install -U pip
```

```
Requirement already satisfied: numpy in c:\users\user\anaconda3\lib\site-packages (1.24.2)
Requirement already satisfied: scipy in c:\users\user\anaconda3\lib\site-packages (1.9.1)
Requirement already satisfied: gymnasium<0.29 in c:\users\user\anaconda3\lib\site-packages (0.28.1)
Requirement already satisfied: jax-jumpy>=1.0.0 in c:\users\user\anaconda3\lib\site-packages (from gymnasium<0.29) (1.0.0)
Requirement already satisfied: importlib-metadata>=4.8.0 in c:\users\user\anaconda3\lib\site-packages (from gymnasium<0.29) (4.11.3)
Requirement already satisfied: farama-notifications>=0.0.1 in c:\users\user\anaconda3\lib\site-packages (from gymnasium<0.29) (0.0.4)
Requirement already satisfied: cloudpickle>=1.2.0 in c:\users\user\anaconda3\lib\site-packages (from gymnasium<0.29) (2.0.0)
Requirement already satisfied: typing-extensions>=4.3.0 in c:\users\user\anaconda3\lib\site-packages (from gymnasium<0.29) (4.3.0)
Requirement already satisfied: zipp>=0.5 in c:\users\user\anaconda3\lib\site-packages (from importlib-metadata>=4.8.0->gymnasium<0.29) (3.8.0)
Requirement already satisfied: torch in c:\users\user\anaconda3\lib\site-packages (2.1.0)
Requirement already satisfied: torchvision in c:\users\user\anaconda3\lib\site-packages (0.16.0)
Requirement already satisfied: filelock in c:\users\user\anaconda3\lib\site-packages (from torch) (3.6.0)
Requirement already satisfied: fsspec in c:\users\user\anaconda3\lib\site-packages (from torch) (2022.7.1)
Requirement already satisfied: networkx in c:\users\user\anaconda3\lib\site-packages (from torch) (2.8.4)
Requirement already satisfied: sympy in c:\users\user\anaconda3\lib\site-packages (from torch) (1.10.1)
Requirement already satisfied: typing-extensions in c:\users\user\anaconda3\lib\site-packages (from torch) (4.3.0)
Requirement already satisfied: jinja2 in c:\users\user\anaconda3\lib\site-packages (from torch) (2.11.3)
Requirement already satisfied: requests in c:\users\user\anaconda3\lib\site-packages (from torchvision) (2.28.1)
Requirement already satisfied: numpy in c:\users\user\anaconda3\lib\site-packages (from torchvision) (1.24.2)
Requirement already satisfied: pillow!=8.3.*,>=5.3.0 in c:\users\user\anaconda3\lib\site-packages (from torchvision) (9.2.0)
Requirement already satisfied: MarkupSafe>=0.23 in c:\users\user\anaconda3\lib\site-packages (from jinja2->torch) (2.0.1)
Requirement already satisfied: charset-normalizer<3,>=2 in c:\users\user\anaconda3\lib\site-packages (from requests->torchvision) (2.0.4)
Requirement already satisfied: urllib3<1.27,>=1.21.1 in c:\users\user\anaconda3\lib\site-packages (from requests->torchvision) (1.26.11)
Requirement already satisfied: certifi>=2017.4.17 in c:\users\user\anaconda3\lib\site-packages (from requests->torchvision) (2022.9.24)
Requirement already satisfied: idna<4,>=2.5 in c:\users\user\anaconda3\lib\site-packages (from requests->torchvision) (3.3)
Requirement already satisfied: mpmath>=0.19 in c:\users\user\anaconda3\lib\site-packages (from sympy->torch) (1.2.1)
Requirement already satisfied: mediapy in c:\users\user\anaconda3\lib\site-packages (1.1.9)
Requirement already satisfied: Pillow in c:\users\user\anaconda3\lib\site-packages (from mediapy) (9.2.0)
Requirement already satisfied: matplotlib in c:\users\user\anaconda3\lib\site-packages (from mediapy) (3.5.2)
Requirement already satisfied: numpy in c:\users\user\anaconda3\lib\site-packages (from mediapy) (1.24.2)
```

```
Requirement already satisfied: ipython in c:\users\user\anaconda3\lib\site-packages  
(from mediapy) (7.31.1)  
Requirement already satisfied: jedi>=0.16 in c:\users\user\anaconda3\lib\site-packages  
(from ipython->mediapy) (0.18.1)  
Requirement already satisfied: backcall in c:\users\user\anaconda3\lib\site-packages  
(from ipython->mediapy) (0.2.0)  
Requirement already satisfied: pickleshare in c:\users\user\anaconda3\lib\site-packages  
(from ipython->mediapy) (0.7.5)  
Requirement already satisfied: pygments in c:\users\user\anaconda3\lib\site-packages  
(from ipython->mediapy) (2.11.2)  
Requirement already satisfied: decorator in c:\users\user\anaconda3\lib\site-packages  
(from ipython->mediapy) (5.1.1)  
Requirement already satisfied: traitlets>=4.2 in c:\users\user\anaconda3\lib\site-packages  
(from ipython->mediapy) (5.1.1)  
Requirement already satisfied: matplotlib-inline in c:\users\user\anaconda3\lib\site-packages  
(from ipython->mediapy) (0.1.6)  
Requirement already satisfied: colorama in c:\users\user\anaconda3\lib\site-packages  
(from ipython->mediapy) (0.4.5)  
Requirement already satisfied: setuptools>=18.5 in c:\users\user\anaconda3\lib\site-packages  
(from ipython->mediapy) (63.4.1)  
Requirement already satisfied: prompt-toolkit!=3.0.0,!=3.0.1,<3.1.0,>=2.0.0 in c:\use  
rs\user\anaconda3\lib\site-packages (from ipython->mediapy) (3.0.20)  
Requirement already satisfied: python-dateutil>=2.7 in c:\users\user\anaconda3\lib\si  
te-packages (from matplotlib->mediapy) (2.8.2)  
Requirement already satisfied: packaging>=20.0 in c:\users\user\anaconda3\lib\site-p  
ackages (from matplotlib->mediapy) (21.3)  
Requirement already satisfied: pyparsing>=2.2.1 in c:\users\user\anaconda3\lib\site-p  
ackages (from matplotlib->mediapy) (3.0.9)  
Requirement already satisfied: cycler>=0.10 in c:\users\user\anaconda3\lib\site-p  
ackages (from matplotlib->mediapy) (0.11.0)  
Requirement already satisfied: fonttools>=4.22.0 in c:\users\user\anaconda3\lib\site-  
packages (from matplotlib->mediapy) (4.25.0)  
Requirement already satisfied: kiwisolver>=1.0.1 in c:\users\user\anaconda3\lib\site-  
packages (from matplotlib->mediapy) (1.4.2)  
Requirement already satisfied: parso<0.9.0,>=0.8.0 in c:\users\user\anaconda3\lib\si  
te-packages (from jedi>=0.16->ipython->mediapy) (0.8.3)  
Requirement already satisfied: wcwidth in c:\users\user\anaconda3\lib\site-packages  
(from prompt-toolkit!=3.0.0,!=3.0.1,<3.1.0,>=2.0.0->ipython->mediapy) (0.2.5)  
Requirement already satisfied: six>=1.5 in c:\users\user\anaconda3\lib\site-packages  
(from python-dateutil>=2.7->matplotlib->mediapy) (1.16.0)
```

```
In [3]: # Install MetaDrive, a lightweight driving simulator
```

```
import sys  
  
if sys.version_info.minor >= 12:  
    raise ValueError("MetaDrive only supports python<3.12.0.")  
  
!pip install "git+https://github.com/metadrive/metadrive"
```

```
Collecting git+https://github.com/metadrive/metadrive
  Cloning https://github.com/metadrive/metadrive to c:\users\user\appdata\local\temp\pip-req-build-iq1uvx4b
    Resolved https://github.com/metadrive/metadrive to commit 81ce062e80e9ce3dd8f1a4692f16eb45d3e8a43c
      Preparing metadata (setup.py): started
      Preparing metadata (setup.py): finished with status 'done'
Requirement already satisfied: requests in c:\users\user\anaconda3\lib\site-packages (from metadrive-simulator==0.4.1.2) (2.28.1)
Requirement already satisfied: gymnasium<0.29,>=0.28 in c:\users\user\anaconda3\lib\site-packages (from metadrive-simulator==0.4.1.2) (0.28.1)
Requirement already satisfied: numpy<=1.24.2,>=1.21.6 in c:\users\user\anaconda3\lib\site-packages (from metadrive-simulator==0.4.1.2) (1.24.2)
Requirement already satisfied: matplotlib in c:\users\user\anaconda3\lib\site-packages (from metadrive-simulator==0.4.1.2) (3.5.2)
Requirement already satisfied: pandas in c:\users\user\anaconda3\lib\site-packages (from metadrive-simulator==0.4.1.2) (1.5.2)
Requirement already satisfied: pygame in c:\users\user\anaconda3\lib\site-packages (from metadrive-simulator==0.4.1.2) (2.5.2)
Requirement already satisfied: tqdm in c:\users\user\anaconda3\lib\site-packages (from metadrive-simulator==0.4.1.2) (4.64.1)
Requirement already satisfied: yapf in c:\users\user\anaconda3\lib\site-packages (from metadrive-simulator==0.4.1.2) (0.31.0)
Requirement already satisfied: seaborn in c:\users\user\anaconda3\lib\site-packages (from metadrive-simulator==0.4.1.2) (0.11.2)
Requirement already satisfied: progressbar in c:\users\user\anaconda3\lib\site-packages (from metadrive-simulator==0.4.1.2) (2.5)
Requirement already satisfied: panda3d==1.10.13 in c:\users\user\anaconda3\lib\site-packages (from metadrive-simulator==0.4.1.2) (1.10.13)
Requirement already satisfied: panda3d-gltf==0.13 in c:\users\user\anaconda3\lib\site-packages (from metadrive-simulator==0.4.1.2) (0.13)
Requirement already satisfied: panda3d-simplepbr in c:\users\user\anaconda3\lib\site-packages (from metadrive-simulator==0.4.1.2) (0.10)
Requirement already satisfied: pillow in c:\users\user\anaconda3\lib\site-packages (from metadrive-simulator==0.4.1.2) (9.2.0)
Requirement already satisfied: pytest in c:\users\user\anaconda3\lib\site-packages (from metadrive-simulator==0.4.1.2) (7.1.2)
Requirement already satisfied: opencv-python in c:\users\user\anaconda3\lib\site-packages (from metadrive-simulator==0.4.1.2) (4.8.1.78)
Requirement already satisfied: lxml in c:\users\user\anaconda3\lib\site-packages (from metadrive-simulator==0.4.1.2) (4.9.1)
Requirement already satisfied: scipy in c:\users\user\anaconda3\lib\site-packages (from metadrive-simulator==0.4.1.2) (1.9.1)
Requirement already satisfied: psutil in c:\users\user\anaconda3\lib\site-packages (from metadrive-simulator==0.4.1.2) (5.9.0)
Requirement already satisfied: geopandas in c:\users\user\anaconda3\lib\site-packages (from metadrive-simulator==0.4.1.2) (0.14.0)
Requirement already satisfied: shapely in c:\users\user\anaconda3\lib\site-packages (from metadrive-simulator==0.4.1.2) (2.0.2)
Requirement already satisfied: filelock in c:\users\user\anaconda3\lib\site-packages (from metadrive-simulator==0.4.1.2) (3.6.0)
Requirement already satisfied: cloudpickle>=1.2.0 in c:\users\user\anaconda3\lib\site-packages (from gymnasium<0.29,>=0.28->metadrive-simulator==0.4.1.2) (2.0.0)
Requirement already satisfied: jax-jumpy>=1.0.0 in c:\users\user\anaconda3\lib\site-packages (from gymnasium<0.29,>=0.28->metadrive-simulator==0.4.1.2) (1.0.0)
Requirement already satisfied: importlib-metadata>=4.8.0 in c:\users\user\anaconda3\lib\site-packages (from gymnasium<0.29,>=0.28->metadrive-simulator==0.4.1.2) (4.11.3)
Requirement already satisfied: farama-notifications>=0.0.1 in c:\users\user\anaconda3\lib\site-packages (from gymnasium<0.29,>=0.28->metadrive-simulator==0.4.1.2) (0.0.4)
Requirement already satisfied: typing-extensions>=4.3.0 in c:\users\user\anaconda3\li
```

```
b\site-packages (from gymnasium<0.29,>=0.28->metadrive-simulator==0.4.1.2) (4.3.0)
Requirement already satisfied: pyproj>=3.3.0 in c:\users\user\anaconda3\lib\site-packages (from geopandas->metadrive-simulator==0.4.1.2) (3.6.1)
Requirement already satisfied: fiona>=1.8.21 in c:\users\user\anaconda3\lib\site-packages (from geopandas->metadrive-simulator==0.4.1.2) (1.9.5)
Requirement already satisfied: packaging in c:\users\user\anaconda3\lib\site-packages (from geopandas->metadrive-simulator==0.4.1.2) (21.3)
Requirement already satisfied: pytz>=2020.1 in c:\users\user\anaconda3\lib\site-packages (from pandas->metadrive-simulator==0.4.1.2) (2022.1)
Requirement already satisfied: python-dateutil>=2.8.1 in c:\users\user\anaconda3\lib\site-packages (from pandas->metadrive-simulator==0.4.1.2) (2.8.2)
Requirement already satisfied: fonttools>=4.22.0 in c:\users\user\anaconda3\lib\site-packages (from matplotlib->metadrive-simulator==0.4.1.2) (4.25.0)
Requirement already satisfied: pyparsing>=2.2.1 in c:\users\user\anaconda3\lib\site-packages (from matplotlib->metadrive-simulator==0.4.1.2) (3.0.9)
Requirement already satisfied: cycler>=0.10 in c:\users\user\anaconda3\lib\site-packages (from matplotlib->metadrive-simulator==0.4.1.2) (0.11.0)
Requirement already satisfied: kiwisolver>=1.0.1 in c:\users\user\anaconda3\lib\site-packages (from matplotlib->metadrive-simulator==0.4.1.2) (1.4.2)
Requirement already satisfied: attrs>=19.2.0 in c:\users\user\anaconda3\lib\site-packages (from pytest->metadrive-simulator==0.4.1.2) (21.4.0)
Requirement already satisfied: configparser in c:\users\user\anaconda3\lib\site-packages (from pytest->metadrive-simulator==0.4.1.2) (1.1.1)
Requirement already satisfied: pluggy<2.0,>=0.12 in c:\users\user\anaconda3\lib\site-packages (from pytest->metadrive-simulator==0.4.1.2) (1.0.0)
Requirement already satisfied: py>=1.8.2 in c:\users\user\anaconda3\lib\site-packages (from pytest->metadrive-simulator==0.4.1.2) (1.11.0)
Requirement already satisfied: tomli>=1.0.0 in c:\users\user\anaconda3\lib\site-packages (from pytest->metadrive-simulator==0.4.1.2) (2.0.1)
Requirement already satisfied: atomicwrites>=1.0 in c:\users\user\anaconda3\lib\site-packages (from pytest->metadrive-simulator==0.4.1.2) (1.4.0)
Requirement already satisfied: colorama in c:\users\user\anaconda3\lib\site-packages (from pytest->metadrive-simulator==0.4.1.2) (0.4.5)
Requirement already satisfied: charset-normalizer<3,>=2 in c:\users\user\anaconda3\lib\site-packages (from requests->metadrive-simulator==0.4.1.2) (2.0.4)
Requirement already satisfied: urllib3<1.27,>=1.21.1 in c:\users\user\anaconda3\lib\site-packages (from requests->metadrive-simulator==0.4.1.2) (1.26.11)
Requirement already satisfied: idna<4,>=2.5 in c:\users\user\anaconda3\lib\site-packages (from requests->metadrive-simulator==0.4.1.2) (3.3)
Requirement already satisfied: certifi>=2017.4.17 in c:\users\user\anaconda3\lib\site-packages (from requests->metadrive-simulator==0.4.1.2) (2022.9.24)
Requirement already satisfied: six in c:\users\user\anaconda3\lib\site-packages (from fiona>=1.8.21->geopandas->metadrive-simulator==0.4.1.2) (1.16.0)
Requirement already satisfied: cligj>=0.5 in c:\users\user\anaconda3\lib\site-packages (from fiona>=1.8.21->geopandas->metadrive-simulator==0.4.1.2) (0.7.2)
Requirement already satisfied: setuptools in c:\users\user\anaconda3\lib\site-packages (from fiona>=1.8.21->geopandas->metadrive-simulator==0.4.1.2) (63.4.1)
Requirement already satisfied: click~=8.0 in c:\users\user\anaconda3\lib\site-packages (from fiona>=1.8.21->geopandas->metadrive-simulator==0.4.1.2) (8.0.4)
Requirement already satisfied: click-plugins>=1.0 in c:\users\user\anaconda3\lib\site-packages (from fiona>=1.8.21->geopandas->metadrive-simulator==0.4.1.2) (1.1.1)
Requirement already satisfied: zipp>=0.5 in c:\users\user\anaconda3\lib\site-packages (from importlib-metadata>=4.8.0->gymnasium<0.29,>=0.28->metadrive-simulator==0.4.1.2) (3.8.0)
```

Running command git clone --filter=blob:none --quiet https://github.com/metadrivers/e/metadrive 'C:\Users\User\AppData\Local\Temp\pip-req-build-iq1uvx4b'

In [4]: # Test whether MetaDrive is properly installed. No error means the test is passed.
!python -m metadrive.examples.profile_metadrive --num-steps 100

Start to profile the efficiency of MetaDrive with 1000 maps and ~4 vehicles!
 Finish 100/100 simulation steps. Time elapse: 0.3306. Average FPS: 580.2896, Average number of vehicles: 5.000
 Total Time Elapse: 0.331, average FPS: 580.290, average number of vehicles: 5.000.

```
[INFO] MetaDrive version: 0.4.1.2
[INFO] Sensors: [lidar: Lidar(50,), side_detector: SideDetector(), lane_line_detector: LaneLineDetector()]
[INFO] Render Mode: none
[INFO] Assets version: 0.4.1.2
[INFO] Episode ended! Scenario Index: 1872 Reason: out_of_road.
```

Section 1: Building abstract class and helper functions

In [5]: # Run this cell without modification

```
# Import some packages that we need to use
import mediapy as media
import gymnasium as gym
import numpy as np
import pandas as pd
import seaborn as sns
from gymnasium.error import Error
from gymnasium import logger
import torch
import torch.nn as nn
from IPython.display import clear_output
import copy
import time
import pygame
import logging

logging.basicConfig(format='[%(levelname)s] %(message)s')
logger = logging.getLogger()
logger.setLevel(logging.INFO)

def wait(sleep=0.2):
    clear_output(wait=True)
    time.sleep(sleep)

def merge_config(new_config, old_config):
    """Merge the user-defined config with default config"""
    config = copy.deepcopy(old_config)
    if new_config is not None:
        config.update(new_config)
    return config

def test_random_policy(policy, env):
    _acts = set()
    for i in range(1000):
        act = policy(0)
        _acts.add(act)
        assert env.action_space.contains(act), "Out of the bound!"
    if len(_acts) != 1:
        print("[HINT] Though we call self.policy 'random policy', " \
```

```

        "we find that generating action randomly at the beginning " \
        "and then fixing it during updating values period lead to better " \
        "performance. Using purely random policy is not even work! " \
        "We encourage you to investigate this issue."
    )

# We register a non-slippery version of FrozenLake environment.
try:
    gym.register(
        id='FrozenLakeNotSlippery-v1',
        entry_point='gymnasium.envs.toy_text:FrozenLakeEnv',
        kwargs={'map_name': '4x4', 'is_slippery': False},
        max_episode_steps=200,
        reward_threshold=0.78, # optimum = .8196
    )
except Error:
    print("The environment is registered already.")

def _render_helper(env, sleep=0.1):
    ret = env.render()
    if sleep:
        wait(sleep=sleep)
    return ret

def animate(img_array, fps=None):
    """A function that can generate GIF file and show in Notebook."""
    media.show_video(img_array, fps=fps)

def evaluate(policy, num_episodes=1, seed=0, env_name='FrozenLake8x8-v1',
            render=None, existing_env=None, max_episode_length=1000,
            sleep=0.0, verbose=False):
    """This function evaluate the given policy and return the mean episode
    reward.
    :param policy: a function whose input is the observation
    :param num_episodes: number of episodes you wish to run
    :param seed: the random seed
    :param env_name: the name of the environment
    :param render: a boolean flag indicating whether to render policy
    :return: the averaged episode reward of the given policy.
    """
    if existing_env is None:
        render_mode = render if render else None
        env = gym.make(env_name, render_mode=render)
    else:
        env = existing_env
    try:
        rewards = []
        frames = []
        succ_rate = []
        if render:
            num_episodes = 1
        for i in range(num_episodes):
            obs, info = env.reset(seed=seed + i)
            act = policy(obs)
            ep_reward = 0
            for step_count in range(max_episode_length):

```

```

obs, reward, terminated, truncated, info = env.step(act)
done = terminated or truncated

act = policy(obs)
ep_reward += reward

if verbose and step_count % 50 == 0:
    print("Evaluating {} / {} episodes. We are in {} / {} steps. Current episode: {} / {}, step count: {} / {}, max episode length: {}, episode reward: {}".
          format(i + 1, num_episodes, step_count + 1, max_episode_length, ep_reward))

if render == "ansi":
    print(_render_helper(env, sleep))
elif render:
    frames.append(_render_helper(env, sleep))
if done:
    break
rewards.append(ep_reward)
if "arrive_dest" in info:
    succ_rate.append(float(info["arrive_dest"]))
if render:
    env.close()
except Exception as e:
    env.close()
    raise e
finally:
    env.close()
eval_dict = {"frames": frames}
if succ_rate:
    eval_dict["success_rate"] = sum(succ_rate) / len(succ_rate)
return np.mean(rewards), eval_dict

```

In [6]: # Run this cell without modification

```

DEFAULT_CONFIG = dict(
    seed=0,
    max_iteration=20000,
    max_episode_length=200,
    evaluate_interval=10,
    evaluate_num_episodes=10,
    learning_rate=0.001,
    gamma=0.8,
    eps=0.3,
    env_name='FrozenLakeNotSlippery-v1'
)

class AbstractTrainer:
    """This is the abstract class for value-based RL trainer. We will inherit
    the specify algorithm's trainer from this abstract class, so that we can
    reuse the codes.
    """

    def __init__(self, config):
        self.config = merge_config(config, DEFAULT_CONFIG)

        # Create the environment
        self.env_name = self.config['env_name']
        self.env = gym.make(self.env_name)

```

```

# Apply the random seed
self.seed = self.config["seed"]
np.random.seed(self.seed)
self.env.reset(seed=self.seed)

# We set self.obs_dim to the number of possible observation
# if observation space is discrete, otherwise the number
# of observation's dimensions. The same to self.act_dim.
if isinstance(self.env.observation_space, gym.spaces.box.Box):
    assert len(self.env.observation_space.shape) == 1
    self.obs_dim = self.env.observation_space.shape[0]
    self.discrete_obs = False
elif isinstance(self.env.observation_space,
               gym.spaces.discrete.Discrete):
    self.obs_dim = self.env.observation_space.n
    self.discrete_obs = True
else:
    raise ValueError("Wrong observation space!")

if isinstance(self.env.action_space, gym.spaces.box.Box):
    assert len(self.env.action_space.shape) == 1
    self.act_dim = self.env.action_space.shape[0]
elif isinstance(self.env.action_space, gym.spaces.discrete.Discrete):
    self.act_dim = self.env.action_space.n
else:
    raise ValueError("Wrong action space! {}".format(self.env.action_space))

self.eps = self.config['eps']

def process_state(self, state):
    """
    Process the raw observation. For example, we can use this function to
    convert the input state (integer) to a one-hot vector.
    """
    return state

def compute_action(self, processed_state, eps=None):
    """
    Compute the action given the processed state.
    """
    raise NotImplementedError(
        "You need to override the Trainer.compute_action() function.")

def evaluate(self, num_episodes=50, *args, **kwargs):
    """
    Use the function you write to evaluate current policy.
    Return the mean episode reward of 50 episodes.
    """
    if "MetaDrive" in self.env_name:
        kwargs["existing_env"] = self.env
    result, eval_infos = evaluate(self.policy, num_episodes, seed=self.seed,
                                  env_name=self.env_name, *args, **kwargs)
    return result, eval_infos

def policy(self, raw_state, eps=0.0):
    """
    A wrapper function takes raw_state as input and output action.
    """
    return self.compute_action(self.process_state(raw_state), eps=eps)

def train(self, iteration=None):
    """
    Conduct one iteration of learning.
    """
    raise NotImplementedError("You need to override the "
                           "Trainer.train() function.")

```

In [7]: # Run this cell without modification

```

def run(trainer_cls, config=None, reward_threshold=None):
    """Run the trainer and report progress, agnostic to the class of trainer
    :param trainer_cls: A trainer class
    :param config: A dict
    :param reward_threshold: the reward threshold to break the training
    :return: The trained trainer and a dataframe containing learning progress
    """
    if config is None:
        config = {}
    trainer = trainer_cls(config)
    config = trainer.config
    start = now = time.time()
    stats = []
    total_steps = 0

    try:
        for i in range(config['max_iteration'] + 1):
            stat = trainer.train(iteration=i)
            stat = stat or {}
            stats.append(stat)
            if "episode_len" in stat:
                total_steps += stat["episode_len"]
            if i % config['evaluate_interval'] == 0 or \
                i == config["max_iteration"]:
                reward, _ = trainer.evaluate(
                    config.get("evaluate_num_episodes", 50),
                    max_episode_length=config.get("max_episode_length", 1000))
                logger.info("Iter {}, {}episodic return is {:.2f}. {}".format(
                    i,
                    "" if total_steps == 0 else "Step {}, ".format(total_steps),
                    reward,
                    {k: round(np.mean(v), 4) for k, v in stat.items()
                     if not np.isnan(v) and k != "frames"
                     }
                )
                if stat else ""
            )
            now = time.time()
        if reward_threshold is not None and reward > reward_threshold:
            logger.info("Iter {}, episodic return {:.3f} is "
                        "greater than reward threshold {}. Congratulations! Now we"
                        "exit the training process.".format(i, reward, reward_threshold))
            break
    except Exception as e:
        print("Error happens during training: ")
        raise e
    finally:
        if hasattr(trainer.env, "close"):
            trainer.env.close()
            print("Environment is closed.")

    return trainer, stats

```

Section 2: Q-Learning

(20/100 points)

Q-learning is an off-policy algorithm who differs from SARSA in the computing of TD error.

Unlike getting the TD error by running policy to get `next_act` a' and compute:

$$r + \gamma Q(s', a') - Q(s, a)$$

as in SARSA, in Q-learning we compute the TD error via:

$$r + \gamma \max_{a'} Q(s', a') - Q(s, a).$$

The reason we call it "off-policy" is that the next-Q value is not computed for the "behavior policy", instead, it is a "virtual policy" that always takes the best action given current Q values.

Section 2.1: Building Q Learning Trainer

```
In [8]: # Solve the TODOs and remove `pass`  
  
# Managing configurations of your experiments is important for your research.  
  
Q_LEARNING_TRAINER_CONFIG = merge_config(dict(  
    eps=0.3,  
) , DEFAULT_CONFIG)  
  
  
class QLearningTrainer(AbstractTrainer):  
    def __init__(self, config=None):  
        config = merge_config(config, Q_LEARNING_TRAINER_CONFIG)  
        super(QLearningTrainer, self).__init__(config=config)  
        self.gamma = self.config["gamma"]  
        self.eps = self.config["eps"]  
        self.max_episode_length = self.config["max_episode_length"]  
        self.learning_rate = self.config["learning_rate"]  
  
        # build the Q table  
        self.table = np.zeros((self.obs_dim, self.act_dim))  
  
    def compute_action(self, obs, eps=None):  
        """Implement epsilon-greedy policy  
  
        It is a function that take an integer (state / observation)  
        as input and return an interger (action).  
        """  
        if eps is None:  
            eps = self.eps  
  
        # TODO: You need to implement the epsilon-greedy policy here.  
        # Implement epsilon-greedy policy  
        p = np.random.random() # returns between 0 and 1  
        if p < eps:  
            action = np.random.choice(self.act_dim)  
        else:  
            action = np.argmax(self.table[obs])  
  
        return action
```

```

def train(self, iteration=None):
    """Conduct one iteration of learning."""
    obs, info = self.env.reset()
    for t in range(self.max_episode_length):
        act = self.compute_action(obs)

        next_obs, reward, terminated, truncated, info = self.env.step(act)
        done = terminated or truncated

        # TODO: compute the TD error, based on the next observation
        td_error = None
        if done: # doesn't have next state
            td_error = reward - self.table[obs][act] # r - Q(s,a)
        else:
            max_next_action = np.argmax(self.table[next_obs]) # argmax(Q(s'))
            td_target = reward + self.gamma * self.table[next_obs][max_next_action]
            td_error = td_target - self.table[obs][act] # td_target - Q(s,a)

        # TODO: compute the new Q value
        # hint: use TD error, self.learning_rate and old Q value
        new_value = self.table[obs][act] + self.learning_rate * td_error # Q(S, A)

        self.table[obs][act] = new_value
        obs = next_obs

        if done:
            break

```

Section 2.2: Use Q Learning to train agent in FrozenLake

In [9]: # Run this cell without modification

```

q_learning_trainer, _ = run(
    trainer_cls=QLearningTrainer,
    config=dict(
        max_iteration=5000,
        evaluate_interval=50,
        evaluate_num_episodes=50,
        env_name='FrozenLakeNotSlippery-v1'
    ),
    reward_threshold=0.99
)

```

```
[INFO] Iter 0, episodic return is 0.00.  
[INFO] Iter 50, episodic return is 0.00.  
[INFO] Iter 100, episodic return is 0.00.  
[INFO] Iter 150, episodic return is 0.00.  
[INFO] Iter 200, episodic return is 0.00.  
[INFO] Iter 250, episodic return is 0.00.  
[INFO] Iter 300, episodic return is 0.00.  
[INFO] Iter 350, episodic return is 0.00.  
[INFO] Iter 400, episodic return is 0.00.  
[INFO] Iter 450, episodic return is 0.00.  
[INFO] Iter 500, episodic return is 0.00.  
[INFO] Iter 550, episodic return is 0.00.  
[INFO] Iter 600, episodic return is 0.00.  
[INFO] Iter 650, episodic return is 0.00.  
[INFO] Iter 700, episodic return is 0.00.  
[INFO] Iter 750, episodic return is 0.00.  
[INFO] Iter 800, episodic return is 0.00.  
[INFO] Iter 850, episodic return is 0.00.  
[INFO] Iter 900, episodic return is 0.00.  
[INFO] Iter 950, episodic return is 0.00.  
[INFO] Iter 1000, episodic return is 0.00.  
[INFO] Iter 1050, episodic return is 0.00.  
[INFO] Iter 1100, episodic return is 0.00.  
[INFO] Iter 1150, episodic return is 0.00.  
[INFO] Iter 1200, episodic return is 0.00.  
[INFO] Iter 1250, episodic return is 0.00.  
[INFO] Iter 1300, episodic return is 0.00.  
[INFO] Iter 1350, episodic return is 0.00.  
[INFO] Iter 1400, episodic return is 0.00.  
[INFO] Iter 1450, episodic return is 0.00.  
[INFO] Iter 1500, episodic return is 0.00.  
[INFO] Iter 1550, episodic return is 0.00.  
[INFO] Iter 1600, episodic return is 0.00.  
[INFO] Iter 1650, episodic return is 0.00.  
[INFO] Iter 1700, episodic return is 0.00.  
[INFO] Iter 1750, episodic return is 0.00.  
[INFO] Iter 1800, episodic return is 0.00.  
[INFO] Iter 1850, episodic return is 0.00.  
[INFO] Iter 1900, episodic return is 0.00.  
[INFO] Iter 1950, episodic return is 0.00.  
[INFO] Iter 2000, episodic return is 0.00.  
[INFO] Iter 2050, episodic return is 0.00.  
[INFO] Iter 2100, episodic return is 0.00.  
[INFO] Iter 2150, episodic return is 1.00.  
[INFO] Iter 2150, episodic return 1.000 is greater than reward threshold 0.99. Congratulation! Now we exit the training process.
```

Environment is closed.

```
In [10]: # Run this cell without modification  
  
# Render the Learned behavior  
, eval_info = evaluate(  
    policy=q_learning_trainer.policy,  
    num_episodes=1,  
    env_name=q_learning_trainer.env_name,  
    render="rgb_array", # Visualize the behavior here in the cell  
    sleep=0.2 # The time interval between two rendering frames  
)  
animate(eval_info["frames"], fps=2)
```



Section 3: Implement Deep Q Learning in Pytorch

(30 / 100 points)

In this section, we will implement a neural network and train it with Deep Q Learning with Pytorch, a powerful deep learning framework.

If you are not familiar with Pytorch, we suggest you to go through pytorch official quickstart tutorials:

1. [quickstart](#)
2. [tutorial on RL](#)

Different from the Q learning in Section 2, we will implement Deep Q Network (DQN) in this section. The main differences are summarized as follows:

DQN requires an experience replay memory to store the transitions. A replay memory is implemented in the following `ExperienceReplayMemory` class. It contains a certain amount of transitions: `(s_t, a_t, r_t, s_t+1, done_t)`. When the memory is full, the earliest transition is discarded and the latest one is stored.

The replay memory increases the sample efficiency (since each transition might be used multiple times) when solving complex task. However, you may find it learn slowly in this assignment since the CartPole-v1 is a relatively easy environment.

DQN has a delayed-updating target network. DQN maintains another neural network called the target network that has identical structure of the Q network. After a certain amount of steps has been taken, the target network copies the parameters of the Q network to itself. The update of the target network will be much less frequent than the update of the Q network, since the Q network is updated in each step.

The target network is used to stabilize the estimation of the TD error. In DQN, the TD error is estimated as:

$$(r_t + \gamma \max_{a_{t+1}} Q^{target}(s_{t+1}, a_{t+1}) - Q(s_t, a_t))$$

The Q value of the next state is estimated by the target network, not the Q network that is being updated. This mechanism can reduce the variance of gradient because the next Q values is not influenced by the update of current Q network.

Section 3.1: Build DQN trainer

In [11]: *# Solve the TODOs and remove `pass`*

```
from collections import deque
import random

class ExperienceReplayMemory:
    """Store and sample the transitions"""

    def __init__(self, capacity):
        # deque is a useful class which acts like a list but only contain
        # finite elements. When adding new element into the deque will make deque full
        # `maxLen` elements, the oldest element (the index 0 element) will be removed.

        # TODO: uncomment next line.
        self.memory = deque(maxlen=capacity)

    def push(self, transition):
        self.memory.append(transition)

    def sample(self, batch_size):
        return random.sample(self.memory, batch_size)

    def __len__(self):
        return len(self.memory)
```

In [12]: *# Solve the TODOs and remove `pass`*

```
class PytorchModel(nn.Module):
    def __init__(self, num_inputs, num_outputs, hidden_units=100):
        super(PytorchModel, self).__init__()

        # TODO: Build a nn.Sequential object as the neural network with two hidden Lay
        #
        # The first hidden Layer takes `num_inputs`-dim vector as input and has `hidde
        # followed by a ReLU activation function.
        #
        # The second hidden Layer takes `hidden_units`-dim vector as input and has `hi
        # followed by a ReLU activation function.
        #
        # The output Layer takes `hidden_units`-dim vector as input and return `num_out
        self.action_value = nn.Sequential(
            nn.Linear(num_inputs, hidden_units),
            nn.ReLU(),
            nn.Linear(hidden_units, hidden_units),
            nn.ReLU(),
            nn.Linear(hidden_units, num_outputs)
        )
```

```

def forward(self, obs):
    return self.action_value(obs)

# Test
test_pytorch_model = PytorchModel(num_inputs=3, num_outputs=7, hidden_units=123)
assert isinstance(test_pytorch_model.action_value, nn.Module)
assert len(test_pytorch_model.state_dict()) == 6
assert test_pytorch_model.state_dict()["action_value.0.weight"].shape == (123, 3)
print("Name of each parameter vectors: ", test_pytorch_model.state_dict().keys())

print("Test passed!")

```

Name of each parameter vectors: odict_keys(['action_value.0.weight', 'action_value.0.bias', 'action_value.2.weight', 'action_value.2.bias', 'action_value.4.weight', 'action_value.4.bias'])
Test passed!

In [13]: # Solve the TODOs and remove `pass`

```

DQN_CONFIG = merge_config(dict(
    parameter_std=0.01,
    learning_rate=0.001,
    hidden_dim=100,
    clip_norm=1.0,
    clip_gradient=True,
    max_iteration=1000,
    max_episode_length=1000,
    evaluate_interval=100,
    gamma=0.99,
    eps=0.3,
    memory_size=50000,
    learn_start=5000,
    batch_size=32,
    target_update_freq=500, # in steps
    learn_freq=1, # in steps
    n=1,
    env_name="CartPole-v1",
), Q_LEARNING_TRAINER_CONFIG)

def to_tensor(x):
    """A helper function to transform a numpy array to a Pytorch Tensor"""
    if isinstance(x, np.ndarray):
        x = torch.from_numpy(x).type(torch.float32)
    assert isinstance(x, torch.Tensor)
    if x.dim() == 3 or x.dim() == 1:
        x = x.unsqueeze(0)
    assert x.dim() == 2 or x.dim() == 4, x.shape
    return x

class DQNTrainer(AbstractTrainer):
    def __init__(self, config):
        config = merge_config(config, DQN_CONFIG)
        self.learning_rate = config["learning_rate"]
        super().__init__(config)

        self.memory = ExperienceReplayMemory(config["memory_size"])

```

```

        self.learn_start = config["learn_start"]
        self.batch_size = config["batch_size"]
        self.target_update_freq = config["target_update_freq"]
        self.clip_norm = config["clip_norm"]
        self.hidden_dim = config["hidden_dim"]
        self.max_episode_length = self.config["max_episode_length"]
        self.learning_rate = self.config["learning_rate"]
        self.gamma = self.config["gamma"]
        self.n = self.config["n"]

        self.step_since_update = 0
        self.total_step = 0

# You need to setup the parameter for your function approximator.
self.initialize_parameters()

def initialize_parameters(self):
    # TODO: Initialize the Q network and the target network using PytorchModel class
    self.network = PytorchModel(self.obs_dim, self.act_dim)
    print("Setting up self.network with obs dim: {} and action dim: {}".format(self.obs_dim, self.act_dim))

    self.network.eval()
    self.network.share_memory()

# Initialize target network to be identical to self.network.
# You should put the weights of self.network into self.target_network.
# TODO: Uncomment next few lines
    self.target_network = PytorchModel(self.obs_dim, self.act_dim)
    self.target_network.load_state_dict(self.network.state_dict())

    self.target_network.eval()

# Build Adam optimizer and MSE Loss.
# TODO: Uncomment next few lines
    self.optimizer = torch.optim.Adam(
        self.network.parameters(), lr=self.learning_rate
    )
    self.loss = nn.MSELoss()

def compute_values(self, processed_state):
    """Compute the value for each potential action. Note that you
    should NOT preprocess the state here."""
    values = self.network(processed_state).detach().numpy()
    return values

def compute_action(self, processed_state, eps=None):
    """Compute the action given the state. Note that the input
    is the processed state."""
    values = self.compute_values(processed_state)
    assert values.ndim == 1, values.shape

    if eps is None:
        eps = self.eps

    if np.random.uniform(0, 1) < eps:
        action = self.env.action_space.sample()
    else:
        action = np.argmax(values)
    return action

```

```
def train(self, iteration=None):
    iteration_string = "" if iteration is None else f"Iter {iteration}: "
    obs, info = self.env.reset()
    processed_obs = self.process_state(obs)
    act = self.compute_action(processed_obs)

    stat = {"loss": [], "success_rate": np.nan}

    for t in range(self.max_episode_length):
        next_obs, reward, terminated, truncated, info = self.env.step(act)
        done = terminated or truncated

        next_processed_obs = self.process_state(next_obs)

        # Push the transition into memory.
        self.memory.push(
            (processed_obs, act, reward, next_processed_obs, done)
        )

        processed_obs = next_processed_obs
        act = self.compute_action(next_processed_obs)
        self.step_since_update += 1
        self.total_step += 1

        if done:
            if "arrive_dest" in info:
                stat["success_rate"] = info["arrive_dest"]
            break

        if t % self.config["learn_freq"] != 0:
            # It's not necessary to update policy in each environmental interaction
            continue

        if len(self.memory) < self.learn_start:
            continue
        elif len(self.memory) == self.learn_start:
            logging.info(
                "{}Current memory contains {} transitions, "
                "start learning!".format(iteration_string, self.learn_start)
            )

    batch = self.memory.sample(self.batch_size)

    # Transform a batch of elements in transitions into tensors.
    state_batch = to_tensor(
        np.stack([transition[0] for transition in batch])
    )
    action_batch = to_tensor(
        np.stack([transition[1] for transition in batch])
    )
    reward_batch = to_tensor(
        np.stack([transition[2] for transition in batch])
    )
    next_state_batch = torch.stack(
        [transition[3] for transition in batch]
    )
    done_batch = to_tensor(
        np.stack([transition[4] for transition in batch])
    )
```

```

with torch.no_grad():

    # TODO: Compute the Q values for the next states by calling target net
    Q_t_plus_one: torch.Tensor = self.target_network(next_state_batch)

    assert isinstance(Q_t_plus_one, torch.Tensor)

    # TODO: Compute the target values for current state.
    # The Q_objective will be used as the objective in the loss function.
    # Hint: Remember to use done_batch.
    Q_objective = reward_batch + (1 - done_batch) * self.gamma * Q_t_plus_
    Q_objective = Q_objective.view(-1)

    assert Q_objective.shape == (self.batch_size,)

    self.network.train() # Set the network to "train" mode.

    # TODO: Collect the Q values in batch.
    # Hint: The network will return the Q values for all actions at a given st
    # So we need to "extract" the Q value for the action we've taken.
    # You need to use torch.gather to manipulate the 2nd dimension of the ret
    # tensor from the network and extract the desired Q values.
    Q_t: torch.Tensor = torch.gather(self.network(state_batch), 1, action_batch)
    Q_t = Q_t.view(-1)

    assert Q_t.shape == Q_objective.shape

    # Update the network
    self.optimizer.zero_grad()
    loss = self.loss(input=Q_t, target=Q_objective)
    stat['loss'].append(loss.item())
    loss.backward()

    # TODO: Apply gradient clipping with pytorch utility. Uncomment next line.
    nn.utils.clip_grad_norm_(self.network.parameters(), self.clip_norm)

    self.optimizer.step()
    self.network.eval()

    if len(self.memory) >= self.learn_start and \
        self.step_since_update > self.target_update_freq:
    self.step_since_update = 0

    # TODO: Copy the weights of self.network to self.target_network.
    self.target_network.load_state_dict(self.network.state_dict())

    self.target_network.eval()

    ret = {"loss": np.mean(stat["loss"]), "episode_len": t}
    if "success_rate" in stat:
        ret["success_rate"] = stat["success_rate"]
    return ret

def process_state(self, state):
    return torch.from_numpy(state).type(torch.float32)

def save(self, loc="model.pt"):
    torch.save(self.network.state_dict(), loc)

```

```
def load(self, loc="model.pt"):
    self.network.load_state_dict(torch.load(loc))
```

Section 3.2: Test DQN trainer

In [14]: *# Run this cell without modification*

```
# Build the test trainer.
test_trainer = DQNTTrainer({})

# Test compute_values
fake_state = test_trainer.env.observation_space.sample()
processed_state = test_trainer.process_state(fake_state)
assert processed_state.shape == (test_trainer.obs_dim,), processed_state.shape
values = test_trainer.compute_values(processed_state)
assert values.shape == (test_trainer.act_dim,), values.shape

test_trainer.train()
print("Now your codes should be bug-free.")

- = run(DQNTTrainer, dict(
    max_iteration=20,
    evaluate_interval=10,
    learn_start=100,
    env_name="CartPole-v1",
))
test_trainer.save("test_trainer.pt")
test_trainer.load("test_trainer.pt")

print("Test passed!")
```

Setting up self.network with obs dim: 4 and action dim: 2

```
C:\Users\User\anaconda3\lib\site-packages\numpy\core\fromnumeric.py:3464: RuntimeWarning: Mean of empty slice.
    return _methods._mean(a, axis=axis, dtype=dtype,
C:\Users\User\anaconda3\lib\site-packages\numpy\core\_methods.py:192: RuntimeWarning: invalid value encountered in scalar divide
    ret = ret.dtype.type(ret / rcount)
```

```
[INFO] Iter 0, Step 9, episodic return is 9.40. {'episode_len': 9.0}
[INFO] Iter 8: Current memory contains 100 transitions, start learning!
```

Now your codes should be bug-free.

Setting up self.network with obs dim: 4 and action dim: 2

```
[INFO] Iter 10, Step 112, episodic return is 9.40. {'loss': 0.178, 'episode_len': 9.0}
[INFO] Iter 20, Step 218, episodic return is 10.60. {'loss': 0.0009, 'episode_len': 8.0}
```

Environment is closed.

Test passed!

Section 3.3: Train DQN agents in CartPole

First, we visualize a random agent in CartPole environment.

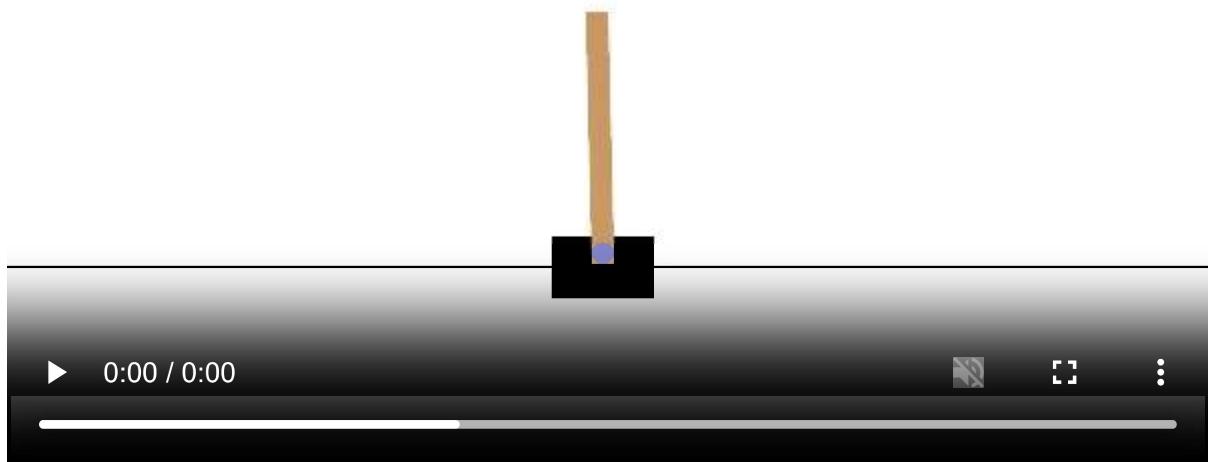
In [15]: *# Run this cell without modification*

```
eval_reward, eval_info = evaluate()
```

```
policy=lambda x: np.random.randint(2),
num_episodes=1,
env_name="CartPole-v1",
render="rgb_array", # Visualize the behavior here in the cell
)

animate(eval_info["frames"])

print("A random agent achieves {} return.".format(eval_reward))
```



A random agent achieves 30.0 return.

In [16]: *# Run this cell without modification*

```
pytorch_trainer, pytorch_stat = run(DQNTrainer, dict(
    max_iteration=5000,
    evaluate_interval=100,
    learning_rate=0.001,
    clip_norm=10.0,
    memory_size=50000,
    learn_start=1000,
    eps=0.1,
    target_update_freq=2000,
    batch_size=128,
    learn_freq=32,
    env_name="CartPole-v1",
), reward_threshold=450.0)

reward, _ = pytorch_trainer.evaluate()
assert reward > 400.0, "Check your codes. " \
                     "Your agent should achieve {} reward in 5000 iterations." \
                     "But it achieve {} reward in evaluation.".format(400.0, reward)

pytorch_trainer.save("dqn_trainer_cartpole.pt")
```

Should solve the task in 10 minutes

```
[INFO] Iter 0, Step 8, episodic return is 9.40. {'episode_len': 8.0}
Setting up self.network with obs dim: 4 and action dim: 2
```

```
[INFO] Iter 100, Step 880, episodic return is 9.40. {'episode_len': 9.0}
[INFO] Iter 200, Step 1784, episodic return is 10.20. {'loss': 0.0014, 'episode_len': 10.0}
[INFO] Iter 300, Step 2699, episodic return is 9.40. {'loss': 0.0817, 'episode_len': 8.0}
[INFO] Iter 400, Step 3608, episodic return is 9.40. {'loss': 0.1022, 'episode_len': 12.0}
[INFO] Iter 500, Step 4512, episodic return is 9.40. {'loss': 0.1372, 'episode_len': 9.0}
[INFO] Iter 600, Step 5377, episodic return is 9.40. {'loss': 0.0529, 'episode_len': 9.0}
[INFO] Iter 700, Step 6327, episodic return is 9.50. {'loss': 0.2304, 'episode_len': 8.0}
[INFO] Iter 800, Step 7404, episodic return is 9.60. {'loss': 0.1802, 'episode_len': 9.0}
[INFO] Iter 900, Step 9566, episodic return is 14.80. {'loss': 0.2251, 'episode_len': 23.0}
[INFO] Iter 1000, Step 13920, episodic return is 114.10. {'loss': 0.3517, 'episode_len': 227.0}
[INFO] Iter 1100, Step 27173, episodic return is 200.80. {'loss': 0.5562, 'episode_len': 234.0}
[INFO] Iter 1200, Step 44581, episodic return is 170.90. {'loss': 1.0156, 'episode_len': 143.0}
[INFO] Iter 1300, Step 59989, episodic return is 158.20. {'loss': 1.6791, 'episode_len': 172.0}
[INFO] Iter 1400, Step 74278, episodic return is 144.10. {'loss': 0.2672, 'episode_len': 155.0}
[INFO] Iter 1500, Step 88183, episodic return is 144.80. {'loss': 0.0639, 'episode_len': 152.0}
[INFO] Iter 1600, Step 102153, episodic return is 142.20. {'loss': 0.0974, 'episode_len': 146.0}
[INFO] Iter 1700, Step 116148, episodic return is 143.70. {'loss': 0.0134, 'episode_len': 129.0}
[INFO] Iter 1800, Step 130211, episodic return is 140.00. {'loss': 0.0156, 'episode_len': 128.0}
[INFO] Iter 1900, Step 143684, episodic return is 137.30. {'loss': 0.0331, 'episode_len': 153.0}
[INFO] Iter 2000, Step 157191, episodic return is 135.80. {'loss': 0.2098, 'episode_len': 150.0}
[INFO] Iter 2100, Step 171066, episodic return is 148.10. {'loss': 0.1178, 'episode_len': 129.0}
[INFO] Iter 2200, Step 185248, episodic return is 140.90. {'loss': 0.0939, 'episode_len': 125.0}
[INFO] Iter 2300, Step 199754, episodic return is 142.00. {'loss': 0.0218, 'episode_len': 156.0}
[INFO] Iter 2400, Step 214286, episodic return is 153.70. {'loss': 0.0133, 'episode_len': 137.0}
[INFO] Iter 2500, Step 229347, episodic return is 146.50. {'loss': 0.0275, 'episode_len': 149.0}
[INFO] Iter 2600, Step 244427, episodic return is 155.10. {'loss': 0.0176, 'episode_len': 175.0}
[INFO] Iter 2700, Step 260090, episodic return is 183.00. {'loss': 0.1468, 'episode_len': 163.0}
[INFO] Iter 2800, Step 276464, episodic return is 179.00. {'loss': 0.0132, 'episode_len': 174.0}
[INFO] Iter 2900, Step 295777, episodic return is 210.30. {'loss': 0.0157, 'episode_len': 190.0}
[INFO] Iter 3000, Step 317918, episodic return is 266.20. {'loss': 0.1372, 'episode_len': 228.0}
[INFO] Iter 3100, Step 344652, episodic return is 171.10. {'loss': 8.9656, 'episode_len': 221.0}
```

```
en': 401.0}
[INFO] Iter 3200, Step 371194, episodic return is 151.70. {'loss': 12.1258, 'episode_len': 166.0}
[INFO] Iter 3300, Step 398342, episodic return is 211.20. {'loss': 3.7496, 'episode_len': 296.0}
[INFO] Iter 3400, Step 427122, episodic return is 305.00. {'loss': 0.0998, 'episode_len': 267.0}
[INFO] Iter 3500, Step 460097, episodic return is 378.00. {'loss': 0.7406, 'episode_len': 282.0}
[INFO] Iter 3600, Step 488518, episodic return is 286.90. {'loss': 0.9118, 'episode_len': 292.0}
[INFO] Iter 3700, Step 519654, episodic return is 334.30. {'loss': 0.6812, 'episode_len': 393.0}
[INFO] Iter 3800, Step 561743, episodic return is 419.20. {'loss': 0.995, 'episode_len': 499.0}
[INFO] Iter 3900, Step 604191, episodic return is 500.00. {'loss': 5.2421, 'episode_len': 499.0}
[INFO] Iter 3900, episodic return 500.000 is greater than reward threshold 450.0. Congratulation! Now we exit the training process.
```

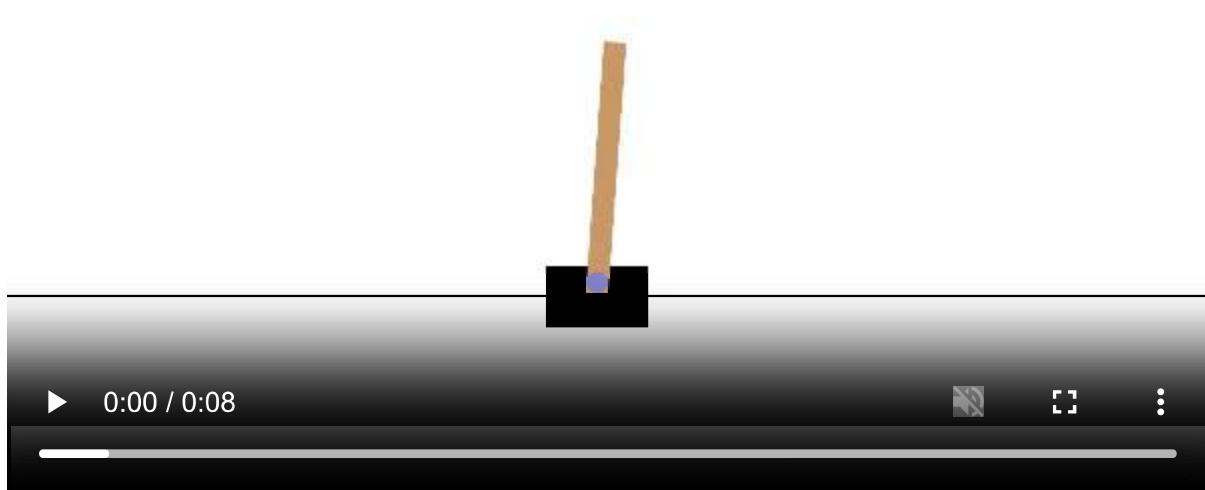
Environment is closed.

In [17]: # Run this cell without modification

```
# Render the Learned behavior
eval_reward, eval_info = evaluate(
    policy=pytorch_trainer.policy,
    num_episodes=1,
    env_name=pytorch_trainer.env_name,
    render="rgb_array", # Visualize the behavior here in the cell
)

animate(eval_info["frames"])

print("DQN agent achieves {}".format(eval_reward))
```



Section 3.4: Train DQN agents in MetaDrive

```
In [18]: # Run this cell without modification

def register_metadrive():
    try:
        from metadrive.envs import MetaDriveEnv
        from metadrive.utils.config import merge_config_with_unknown_keys
    except ImportError as e:
        print("Please install MetaDrive through: pip install git+https://github.com/de"
              "raise e

    env_names = []
    try:
        class MetaDriveEnvTut(gym.Wrapper):
            def __init__(self, config, *args, render_mode=None, **kwargs):
                # Ignore render_mode
                self._render_mode = render_mode
                super().__init__(MetaDriveEnv(config))
                self.action_space = gym.spaces.Discrete(int(np.prod(self.env.action_sp

            def reset(self, *args, seed=None, render_mode=None, options=None, **kwargs):
                # Ignore seed and render_mode
                return self.env.reset(*args, **kwargs)

            def render(self):
                return self.env.render(mode=self._render_mode)

            def _make_env(*args, **kwargs):
                return MetaDriveEnvTut(*args, **kwargs)

        env_name = "MetaDrive-Tut-Easy-v0"
        gym.register(id=env_name, entry_point=_make_env, kwargs={"config": dict(
            map="S",
            start_seed=0,
            num_scenarios=1,
            horizon=200,
            discrete_action=True,
            discrete_steering_dim=3,
            discrete_throttle_dim=3
        )})
        env_names.append(env_name)

        env_name = "MetaDrive-Tut-Hard-v0"
        gym.register(id=env_name, entry_point=_make_env, kwargs={"config": dict(
            map="CCC",
            start_seed=0,
            num_scenarios=10,
            discrete_action=True,
            discrete_steering_dim=5,
            discrete_throttle_dim=5
        )})
        env_names.append(env_name)
    except gym.error.Error as e:
        print("Information when registering MetaDrive: ", e)
    else:
        print("Successfully registered MetaDrive environments: ", env_names)
```

```
In [19]: # Run this cell without modification
register_metadrive()

Successfully registered MetaDrive environments: ['MetaDrive-Tut-Easy-v0', 'MetaDrive-Tut-Hard-v0']
```

```
In [20]: # Run this cell without modification

# Build the test trainer.
test_trainer = DQNTrainer(dict(env_name="MetaDrive-Tut-Easy-v0"))

# Test compute_values
for _ in range(10):
    fake_state = test_trainer.env.observation_space.sample()
    processed_state = test_trainer.process_state(fake_state)
    assert processed_state.shape == (test_trainer.obs_dim,), processed_state.shape
    values = test_trainer.compute_values(processed_state)
    assert values.shape == (test_trainer.act_dim,), values.shape

    test_trainer.train()

print("Now your codes should be bug-free.")
test_trainer.env.close()
del test_trainer
```

```
[INFO] MetaDrive version: 0.4.1.2
[INFO] Sensors: [lidar: Lidar(50,), side_detector: SideDetector(), lane_line_detector: LaneLineDetector()]
[INFO] Render Mode: none
[INFO] Assets version: 0.4.1.2
```

Setting up self.network with obs dim: 259 and action dim: 9

```
[INFO] Episode ended! Scenario Index: 0 Reason: max step
[INFO] Episode ended! Scenario Index: 0 Reason: max step
[INFO] Episode ended! Scenario Index: 0 Reason: max step
[INFO] Episode ended! Scenario Index: 0 Reason: max step
[INFO] Episode ended! Scenario Index: 0 Reason: max step
[INFO] Episode ended! Scenario Index: 0 Reason: max step
[INFO] Episode ended! Scenario Index: 0 Reason: max step
[INFO] Episode ended! Scenario Index: 0 Reason: max step
[INFO] Episode ended! Scenario Index: 0 Reason: max step
[INFO] Episode ended! Scenario Index: 0 Reason: max step
```

Now your codes should be bug-free.

```
In [21]: # Run this cell without modification

env_name = "MetaDrive-Tut-Easy-v0"

pytorch_trainer2, _ = run(DQNTrainer, dict(
    max_episode_length=200,
    max_iteration=5000,
    evaluate_interval=10,
    evaluate_num_episodes=10,
    learning_rate=0.0001,
    clip_norm=10.0,
    memory_size=1000000,
    learn_start=2000,
    eps=0.1,
    target_update_freq=5000,
```

```
learn_freq=16,
batch_size=256,
env_name=env_name
), reward_threshold=120)

pytorch_trainer2.save("dqn_trainer_metadrive_easy.pt")

# Run this cell without modification

# Render the Learned behavior
# NOTE: The Learned agent is marked by green color.
eval_reward, eval_info = evaluate(
    policy=pytorch_trainer2.policy,
    num_episodes=1,
    env_name=pytorch_trainer2.env_name,
    render="topdown", # Visualize the behaviors in top-down view
    verbose=True
)

frames = [pygame.surfarray.array3d(f).swapaxes(0, 1) for f in eval_info["frames"]]

animate(frames)

print("DQN agent achieves {} return in MetaDrive easy environment.".format(eval_reward))
```

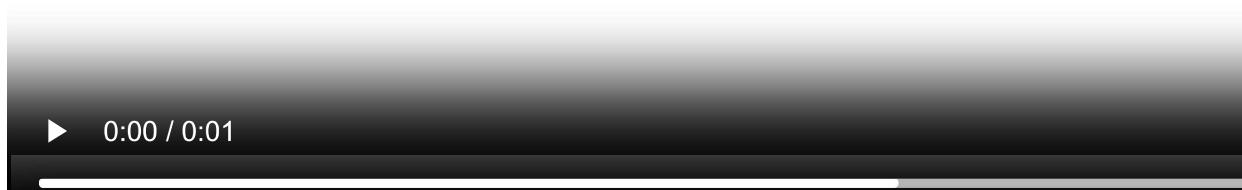
Setting up self.network with obs dim: 259 and action dim: 9

```
[INFO] Iter 0, Step 13, episodic return is -4.48. {'episode_len': 13.0, 'success_rate': 0.0}
[INFO] Iter 10, Step 107, episodic return is -4.48. {'episode_len': 8.0, 'success_rate': 0.0}
[INFO] Iter 20, Step 206, episodic return is -4.48. {'episode_len': 12.0, 'success_rate': 0.0}
[INFO] Iter 30, Step 299, episodic return is -4.48. {'episode_len': 9.0, 'success_rate': 0.0}
[INFO] Iter 40, Step 392, episodic return is -4.48. {'episode_len': 8.0, 'success_rate': 0.0}
[INFO] Iter 50, Step 480, episodic return is -4.48. {'episode_len': 8.0, 'success_rate': 0.0}
[INFO] Iter 60, Step 570, episodic return is -4.48. {'episode_len': 8.0, 'success_rate': 0.0}
[INFO] Iter 70, Step 663, episodic return is -4.48. {'episode_len': 10.0, 'success_rate': 0.0}
[INFO] Iter 80, Step 755, episodic return is -4.48. {'episode_len': 16.0, 'success_rate': 0.0}
[INFO] Iter 90, Step 847, episodic return is -4.48. {'episode_len': 14.0, 'success_rate': 0.0}
[INFO] Iter 100, Step 951, episodic return is -4.48. {'episode_len': 9.0, 'success_rate': 0.0}
[INFO] Iter 110, Step 1046, episodic return is -4.48. {'episode_len': 9.0, 'success_rate': 0.0}
[INFO] Iter 120, Step 1130, episodic return is -4.48. {'episode_len': 9.0, 'success_rate': 0.0}
[INFO] Iter 130, Step 1222, episodic return is -4.48. {'episode_len': 8.0, 'success_rate': 0.0}
[INFO] Iter 140, Step 1314, episodic return is -4.48. {'episode_len': 8.0, 'success_rate': 0.0}
[INFO] Iter 150, Step 1418, episodic return is -4.48. {'episode_len': 10.0, 'success_rate': 0.0}
[INFO] Iter 160, Step 1511, episodic return is -4.48. {'episode_len': 8.0, 'success_rate': 0.0}
[INFO] Iter 170, Step 1611, episodic return is -4.48. {'episode_len': 12.0, 'success_rate': 0.0}
[INFO] Iter 180, Step 1697, episodic return is -4.48. {'episode_len': 9.0, 'success_rate': 0.0}
[INFO] Iter 190, Step 1785, episodic return is -4.48. {'episode_len': 10.0, 'success_rate': 0.0}
[INFO] Iter 200, Step 3021, episodic return is 0.01. {'loss': 1.7653, 'episode_len': 199.0}
[INFO] Iter 210, Step 3987, episodic return is 125.54. {'loss': 0.9338, 'episode_len': 73.0, 'success_rate': 0.0}
[INFO] Iter 210, episodic return 125.539 is greater than reward threshold 120. Congratulation! Now we exit the training process.
```

Environment is closed.

Evaluating 1/1 episodes. We are in 1/1000 steps. Current episode reward: 0.000

Evaluating 1/1 episodes. We are in 51/1000 steps. Current episode reward: 35.980



Section 4: Policy gradient methods - REINFORCE

(30 / 100 points)

Unlike the supervised learning, in RL the optimization objective, the episodic return, is not differentiable w.r.t. the neural network parameters. This can be solved via **Policy Gradient**. It can be proved that policy gradient is an unbiased estimator of the gradient of the objective.

Concretely, let's consider such optimization objective:

$$Q = \mathbb{E}_{\text{possible trajectories}} \sum_t r(a_t, s_t) = \sum_{s_0, a_0, \dots} p(s_0, a_0, \dots, s_t, a_t) r(s_0, a_0, \dots, s_t, a_t) = \sum_{\tau} p(\tau) r(\tau)$$

wherein $\sum_t r(a_t, s_t) = r(\tau)$ is the return of trajectory $\tau = (s_0, a_0, \dots)$. We remove the discount factor for simplicity. Since we want to maximize Q , we can simply compute the gradient of Q w.r.t. parameter θ (which is implicitly included in $p(\tau)$):

$$\nabla_{\theta} Q = \nabla_{\theta} \sum_{\tau} p(\tau) r(\tau) = \sum_{\tau} r(\tau) \nabla_{\theta} p(\tau)$$

wherein we've applied a famous trick: $\nabla_{\theta} p(\tau) = p(\tau) \frac{\nabla_{\theta} p(\tau)}{p(\tau)} = p(\tau) \nabla_{\theta} \log p(\tau)$. Here the $r(\tau)$ will be determined when τ is determined. So it has nothing to do with the policy. We can move it out from the gradient.

Introducing a log term can change the product of probabilities to sum of log probabilities. Now we can expand the log of product above to sum of log:

$$p_{\theta}(\tau) = p(s_0, a_0, \dots) = p(s_0) \prod_t \pi_{\theta}(a_t | s_t) p(s_{t+1} | s_t, a_t)$$

$$\log p_{\theta}(\tau) = \log p(s_0) + \sum_t \log \pi_{\theta}(a_t | s_t) + \sum_t \log p(s_{t+1} | s_t, a_t)$$

You can find that the first and third term are not correlated to the parameter of policy $\pi_{\theta}(\cdot)$. So when we compute $\nabla_{\theta} Q$, we find

$$\nabla_{\theta} Q = \sum_{\tau} r(\tau) \nabla_{\theta} p(\tau) = \sum_{\tau} r(\tau) p(\tau) \nabla_{\theta} \log p(\tau) = \sum_{\tau} p_{\theta}(\tau) \left(\sum_t \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right) r(\tau) d\tau$$

When we sample sufficient amount of data from the environment, the above equation can be estimated via:

$$\nabla_{\theta} Q = \frac{1}{N} \sum_{i=1}^N \left[\left(\sum_t \nabla_{\theta} \log \pi_{\theta}(a_{i,t} | s_{i,t}) \right) \left(\sum_{t'=t}^N \gamma^{t'-t} r(s_{i,t'}, a_{i,t'}) \right) \right]$$

This algorithm is called REINFORCE algorithm, which is a Monte Carlo Policy Gradient algorithm with long history. In this section, we will implement the it using pytorch.

The policy network is composed by two parts:

1. A basic neural network serves as the function approximator. It output raw values parameterizing the action distribution given current observation. We will reuse PytorchModel here.
2. A distribution layer builds upon the neural network to wrap the raw logits output from neural network to a distribution and provides API for sampling action and computing log probability.

Section 4.1: Build REINFORCE

```
In [49]: # Solve the TODOs and remove `pass`  
  
class PGNetwork(nn.Module):  
    def __init__(self, obs_dim, act_dim, hidden_units=128):  
        super(PGNetwork, self).__init__()  
        self.network = PytorchModel(obs_dim, act_dim, hidden_units)  
  
    def forward(self, obs):  
        logit = self.network(obs)  
  
        # TODO: Create an object of the class "torch.distributions.Categorical"  
        # Then sample an action from it.  
        action_distribution = torch.distributions.Categorical(logits=logit)  
        action = action_distribution.sample()  
  
        return action  
  
    def log_prob(self, obs, act):  
        logits = self.network(obs)  
  
        # TODO: Create an object of the class "torch.distributions.Categorical"  
        # Then get the log probability of the action `act` in this distribution.  
        action_distribution = torch.distributions.Categorical(logits=logits)  
        log_prob = action_distribution.log_prob(act)  
  
        return log_prob  
  
# Note that we do not implement GaussianPolicy here. So we can't  
# apply our algorithm to the environment with continuous action.
```

```
In [61]: # Solve the TODOs and remove `pass`
```

```
PG_DEFAULT_CONFIG = merge_config(dict(  
    normalize_advantage=True,  
  
    clip_norm=10.0,  
    clip_gradient=True,  
  
    hidden_units=100,
```

```
max_iteration=1000,  
  
train_batch_size=1000,  
gamma=0.99,  
learning_rate=0.001,  
  
env_name="CartPole-v1",  
  
, DEFAULT_CONFIG)  
  
  
class PGTrainer(AbstractTrainer):  
    def __init__(self, config=None):  
        config = merge_config(config, PG_DEFAULT_CONFIG)  
        super().__init__(config)  
  
        self.iteration = 0  
        self.start_time = time.time()  
        self.iteration_time = self.start_time  
        self.total_timesteps = 0  
        self.total_episodes = 0  
  
        # build the model  
        self.initialize_parameters()  
  
    def initialize_parameters(self):  
        """Build the policy network and related optimizer"""  
        # Detect whether you have GPU or not. Remember to call X.to(self.device)  
        # if necessary.  
        self.device = torch.device(  
            "cuda" if torch.cuda.is_available() else "cpu"  
        )  
  
        # TODO Build the policy network using CategoricalPolicy  
        # Hint: Remember to pass config["hidden_units"], and set policy network  
        # to the device you are using.  
        self.network = PGNetwork(  
            self.obs_dim, self.act_dim,  
            hidden_units=self.config["hidden_units"]  
        ).to(self.device)  
  
        # Build the Adam optimizer.  
        self.optimizer = torch.optim.Adam(  
            self.network.parameters(),  
            lr=self.config["learning_rate"]  
        )  
  
    def to_tensor(self, array):  
        """Transform a numpy array to a pytorch tensor"""  
        return torch.from_numpy(array).type(torch.float32).to(self.device)  
  
    def to_array(self, tensor):  
        """Transform a pytorch tensor to a numpy array"""  
        ret = tensor.cpu().detach().numpy()  
        if ret.size == 1:  
            ret = ret.item()  
        return ret  
  
    def save(self, loc="model.pt"):
```

```

        torch.save(self.network.state_dict(), loc)

    def load(self, loc="model.pt"):
        self.network.load_state_dict(torch.load(loc))

    def compute_action(self, observation, eps=None):
        """Compute the action for single observation. eps is useless here."""
        assert observation.ndim == 1
        # TODO: Sample an action from the action distribution given by the policy.
        # Hint: The input of policy network is a tensor with the first dimension to the
        # batch dimension. Therefore you need to expand the first dimension of the ob
        # and convert it to a tensor before feeding it to the policy network.

        # Convert the observation to a PyTorch tensor and add batch dimension
        obs_tensor = self.to_tensor(observation).unsqueeze(0)
        # action_logits = self.network(obs_tensor)
        # action_Logits = action_Logits.view(-1)

        # Pass the observation through the policy network to get the action distribution
        # action_distribution = torch.distributions.Categorical(Logits=action_Logits)

        # Sample an action from the distribution
        # action = action_distribution.sample().item()
        action_tensor = self.network.forward(obs_tensor)
        action = action_tensor.item()

        # print(f"[DEBUG] action = {action}")

        return action

    def compute_log_probs(self, observation, action):
        """Compute the log probabilities of a batch of state-action pair"""
        # TODO: Use the function of the policy network to get log probs.
        # Hint: Remember to transform the data into tensor before feeding it into the
        # Convert the observation to a PyTorch tensor and add batch dimension
        obs_tensor = self.to_tensor(observation).unsqueeze(0)
        # action_logits = self.network(obs_tensor)
        # action_Logits = action_Logits.view(-1)

        # Pass the observation through the policy network to get the action distribution
        # action_distribution = torch.distributions.Categorical(Logits=action_Logits)

        # Convert the action to a PyTorch tensor
        action_tensor = self.to_tensor(action)

        # Calculate the Log probability of the given action
        log_probs = action_distribution.log_prob(action_tensor)
        log_probs.requires_grad = True

        log_probs = self.network.log_prob(obs_tensor, action_tensor)
        log_probs = log_probs.view(-1)

        # print(f"[DEBUG] log_probs = {log_probs}")

        return log_probs

    def update_network(self, processed_samples):
        """Update the policy network"""
        advantages = self.to_tensor(processed_samples["advantages"])
        flat_obs = np.concatenate(processed_samples["obs"])

```

```

flat_act = np.concatenate(processed_samples["act"])

self.network.train()
self.optimizer.zero_grad()

log_probs = self.compute_log_probs(flat_obs, flat_act)

assert log_probs.shape == advantages.shape, "log_probs shape {} is not " \
                                               "compatible with advantages {}".format(log_probs.shape, advantages.shape)

# TODO: Compute the policy gradient loss.
loss = -torch.mean(log_probs * advantages)
print(f"[DEBUG] Loss = {loss}")

loss.backward()

# Clip the gradient
torch.nn.utils.clip_grad_norm_(
    self.network.parameters(), self.config["clip_gradient"]
)

self.optimizer.step()
self.network.eval()

update_info = {
    "policy_loss": loss.item(),
    "mean_log_prob": torch.mean(log_probs).item(),
    "mean_advantage": torch.mean(advantages).item()
}
return update_info

# ===== Training-related functions =====
def collect_samples(self):
    """Here we define the pipeline to collect sample even though
    any specify functions are not implemented yet.
    """

    iter_timesteps = 0
    iter_episodes = 0
    episode_lens = []
    episode_rewards = []
    episode_obs_list = []
    episode_act_list = []
    episode_reward_list = []
    success_list = []

    while iter_timesteps <= self.config["train_batch_size"]:
        obs_list, act_list, reward_list = [], [], []
        obs, info = self.env.reset()
        steps = 0
        episode_reward = 0
        while True:
            act = self.compute_action(obs)

            next_obs, reward, terminated, truncated, step_info = self.env.step(act)
            done = terminated or truncated

            obs_list.append(obs)
            act_list.append(act)
            reward_list.append(reward)

```

```

        obs = next_obs.copy()
        steps += 1
        episode_reward += reward
        if done or steps > self.config["max_episode_length"]:
            if "arrive_dest" in step_info:
                success_list.append(step_info["arrive_dest"])
            break
        iter_timesteps += steps
        iter_episodes += 1
        episode_rewards.append(episode_reward)
        episode_lens.append(steps)
        episode_obs_list.append(np.array(obs_list, dtype=np.float32))
        episode_act_list.append(np.array(act_list, dtype=np.float32))
        episode_reward_list.append(np.array(reward_list, dtype=np.float32))

    # The return `samples` is a dict that contains several key-value pair.
    # The value of each key-value pair is a list storing the data in one episode.
    samples = {
        "obs": episode_obs_list,
        "act": episode_act_list,
        "reward": episode_reward_list
    }

    sample_info = {
        "iter_timesteps": iter_timesteps,
        "iter_episodes": iter_episodes,
        "performance": np.mean(episode_rewards), # help drawing figures
        "ep_len": float(np.mean(episode_lens)),
        "ep_ret": float(np.mean(episode_rewards)),
        "episode_len": sum(episode_lens),
        "success_rate": np.mean(success_list)
    }
    return samples, sample_info

def process_samples(self, samples):
    """Process samples and add advantages in it"""
    values = []
    for reward_list in samples["reward"]:
        # reward_list contains rewards in one episode
        returns = np.zeros_like(reward_list, dtype=np.float32)
        Q = 0

        # TODO: Scan the reward_list in a reverse order and compute the
        # discounted return at each time step. Fill the array `returns`
        for t in range(len(reward_list) - 1, -1, -1):
            Q = Q * self.config["gamma"] + reward_list[t]
            returns[t] = Q

        values.append(returns)

    # We call the values advantage here.
    advantages = np.concatenate(values)

    if self.config["normalize_advantage"]:
        # TODO: normalize the advantage so that it's mean is
        # almost 0 and the its standard deviation is almost 1.
        mean_advantage = np.mean(advantages)
        std_advantage = np.std(advantages)
        if std_advantage == 0:
            advantages = (advantages - mean_advantage)

```

```

        else:
            advantages = (advantages - mean_advantage) / std_advantage

        samples["advantages"] = advantages
        return samples, {}

# ===== Training iteration =====
def train(self, iteration=None):
    """Here we defined the training pipeline using the abstract
    functions."""
    info = dict(iteration=iteration)

    # Collect samples
    samples, sample_info = self.collect_samples()
    info.update(sample_info)

    # Process samples
    processed_samples, processed_info = self.process_samples(samples)
    info.update(processed_info)

    # Update the model
    update_info = self.update_network(processed_samples)
    info.update(update_info)

    now = time.time()
    self.iteration += 1
    self.total_timesteps += info.pop("iter_timesteps")
    self.total_episodes += info.pop("iter_episodes")

    # info["iter_time"] = now - self.iteration_time
    # info["total_time"] = now - self.start_time
    info["total_episodes"] = self.total_episodes
    info["total_timesteps"] = self.total_timesteps
    self.iteration_time = now

    # print("INFO: ", info)

    return info

```

Section 4.2: Test REINFORCE

In [62]:

```
# Run this cell without modification

# Test advantage computing
test_trainer = PGTrainer({"normalize_advantage": False})
test_trainer.train()
fake_sample = {"reward": [[2, 2, 2, 2, 2]]}
np.testing.assert_almost_equal(
    test_trainer.process_samples(fake_sample)[0]["reward"][0],
    fake_sample["reward"][0]
)
np.testing.assert_almost_equal(
    test_trainer.process_samples(fake_sample)[0]["advantages"],
    np.array([9.80199, 7.880798, 5.9402, 3.98, 2.], dtype=np.float32)
)

# Test advantage normalization
test_trainer = PGTrainer(

```

```

    {"normalize_advantage": True, "env_name": "CartPole-v1"})
test_adv = test_trainer.process_samples(fake_sample)[0]["advantages"]
np.testing.assert_almost_equal(test_adv.mean(), 0.0)
np.testing.assert_almost_equal(test_adv.std(), 1.0)

# Test the shape of functions' returns
fake_observation = np.array([
    test_trainer.env.observation_space.sample() for i in range(10)
])
fake_action = np.array([
    test_trainer.env.action_space.sample() for i in range(10)
])
assert test_trainer.to_tensor(fake_observation).shape == torch.Size([10, 4])
assert np.array(test_trainer.compute_action(fake_observation[0])).shape == ()
assert test_trainer.compute_log_probs(fake_observation, fake_action).shape == \
    torch.Size([10])

print("Test Passed!")

```

Test Passed!

Section 4.3: Train REINFORCE in CartPole and see the impact of advantage normalization

In [63]: # Run this cell without modification

```

pg_trainer_no_na, pg_result_no_na = run(PGTrainer, dict(
    learning_rate=0.001,
    max_episode_length=200,
    train_batch_size=200,
    env_name="CartPole-v1",
    normalize_advantage=False, # <== Here!

    evaluate_interval=10,
    evaluate_num_episodes=10,
), 195.0)

```

```
[INFO] Iter 0, Step 211, episodic return is 20.60. {'iteration': 0.0, 'performance': 21.1, 'ep_len': 21.1, 'ep_ret': 21.1, 'episode_len': 211.0, 'policy_loss': 8.3189, 'mean_log_prob': -0.6923, 'mean_advantage': 11.9656, 'total_episodes': 10.0, 'total_timesteps': 211.0}
[INFO] Iter 10, Step 2445, episodic return is 21.90. {'iteration': 10.0, 'performance': 35.5, 'ep_len': 35.5, 'ep_ret': 35.5, 'episode_len': 213.0, 'policy_loss': 12.575, 'mean_log_prob': -0.6942, 'mean_advantage': 18.6278, 'total_episodes': 84.0, 'total_timesteps': 2445.0}
[INFO] Iter 20, Step 4648, episodic return is 66.30. {'iteration': 20.0, 'performance': 56.5, 'ep_len': 56.5, 'ep_ret': 56.5, 'episode_len': 226.0, 'policy_loss': 18.4315, 'mean_log_prob': -0.6334, 'mean_advantage': 29.3347, 'total_episodes': 144.0, 'total_timesteps': 4648.0}
[INFO] Iter 30, Step 7047, episodic return is 34.40. {'iteration': 30.0, 'performance': 40.6667, 'ep_len': 40.6667, 'ep_ret': 40.6667, 'episode_len': 244.0, 'policy_loss': 12.1378, 'mean_log_prob': -0.619, 'mean_advantage': 19.4421, 'total_episodes': 198.0, 'total_timesteps': 7047.0}
[INFO] Iter 40, Step 9511, episodic return is 50.90. {'iteration': 40.0, 'performance': 65.5, 'ep_len': 65.5, 'ep_ret': 65.5, 'episode_len': 262.0, 'policy_loss': 18.6666, 'mean_log_prob': -0.6211, 'mean_advantage': 29.3126, 'total_episodes': 244.0, 'total_timesteps': 9511.0}
[INFO] Iter 50, Step 11854, episodic return is 93.10. {'iteration': 50.0, 'performance': 61.75, 'ep_len': 61.75, 'ep_ret': 61.75, 'episode_len': 247.0, 'policy_loss': 17.51, 'mean_log_prob': -0.6359, 'mean_advantage': 27.8136, 'total_episodes': 277.0, 'total_timesteps': 11854.0}
[INFO] Iter 60, Step 14208, episodic return is 74.50. {'iteration': 60.0, 'performance': 52.0, 'ep_len': 52.0, 'ep_ret': 52.0, 'episode_len': 208.0, 'policy_loss': 14.3506, 'mean_log_prob': -0.6049, 'mean_advantage': 23.5017, 'total_episodes': 314.0, 'total_timesteps': 14208.0}
[INFO] Iter 70, Step 16723, episodic return is 146.00. {'iteration': 70.0, 'performance': 137.0, 'ep_len': 137.0, 'ep_ret': 137.0, 'episode_len': 274.0, 'policy_loss': 26.2932, 'mean_log_prob': -0.5712, 'mean_advantage': 46.2098, 'total_episodes': 344.0, 'total_timesteps': 16723.0}
[INFO] Iter 80, Step 19384, episodic return is 133.50. {'iteration': 80.0, 'performance': 87.0, 'ep_len': 87.0, 'ep_ret': 87.0, 'episode_len': 261.0, 'policy_loss': 20.8993, 'mean_log_prob': -0.6073, 'mean_advantage': 33.8355, 'total_episodes': 368.0, 'total_timesteps': 19384.0}
[INFO] Iter 90, Step 21983, episodic return is 128.00. {'iteration': 90.0, 'performance': 105.5, 'ep_len': 105.5, 'ep_ret': 105.5, 'episode_len': 211.0, 'policy_loss': 24.6451, 'mean_log_prob': -0.5942, 'mean_advantage': 42.2764, 'total_episodes': 388.0, 'total_timesteps': 21983.0}
[INFO] Iter 100, Step 25112, episodic return is 143.00. {'iteration': 100.0, 'performance': 95.3333, 'ep_len': 95.3333, 'ep_ret': 95.3333, 'episode_len': 286.0, 'policy_loss': 24.9487, 'mean_log_prob': -0.6024, 'mean_advantage': 41.1685, 'total_episodes': 412.0, 'total_timesteps': 25112.0}
[INFO] Iter 110, Step 27489, episodic return is 126.10. {'iteration': 110.0, 'performance': 136.0, 'ep_len': 136.0, 'ep_ret': 136.0, 'episode_len': 272.0, 'policy_loss': 28.8475, 'mean_log_prob': -0.6011, 'mean_advantage': 49.8641, 'total_episodes': 429.0, 'total_timesteps': 27489.0}
[INFO] Iter 120, Step 29897, episodic return is 145.80. {'iteration': 120.0, 'performance': 108.0, 'ep_len': 108.0, 'ep_ret': 108.0, 'episode_len': 324.0, 'policy_loss': 26.6677, 'mean_log_prob': -0.589, 'mean_advantage': 46.3074, 'total_episodes': 446.0, 'total_timesteps': 29897.0}
[INFO] Iter 130, Step 32050, episodic return is 169.30. {'iteration': 130.0, 'performance': 201.0, 'ep_len': 201.0, 'ep_ret': 201.0, 'episode_len': 201.0, 'policy_loss': 32.1502, 'mean_log_prob': -0.5619, 'mean_advantage': 57.2793, 'total_episodes': 461.0, 'total_timesteps': 32050.0}
[INFO] Iter 140, Step 34591, episodic return is 178.60. {'iteration': 140.0, 'performance': 138.5, 'ep_len': 138.5, 'ep_ret': 138.5, 'episode_len': 277.0, 'policy_loss': 27.843, 'mean_log_prob': -0.5574, 'mean_advantage': 49.911, 'total_episodes': 477.0, 'total_timesteps': 34591.0}
```

```
[INFO] Iter 150, Step 37134, episodic return is 186.60. {'iteration': 150.0, 'performance': 199.0, 'ep_len': 199.0, 'ep_ret': 199.0, 'episode_len': 398.0, 'policy_loss': 32.4533, 'mean_log_prob': -0.5653, 'mean_advantage': 56.9853, 'total_episodes': 493.0, 'total_timesteps': 37134.0}
[INFO] Iter 160, Step 39813, episodic return is 138.50. {'iteration': 160.0, 'performance': 192.5, 'ep_len': 192.5, 'ep_ret': 192.5, 'episode_len': 385.0, 'policy_loss': 31.9613, 'mean_log_prob': -0.5694, 'mean_advantage': 56.0284, 'total_episodes': 508.0, 'total_timesteps': 39813.0}
[INFO] Iter 170, Step 42566, episodic return is 159.50. {'iteration': 170.0, 'performance': 165.0, 'ep_len': 165.0, 'ep_ret': 165.0, 'episode_len': 330.0, 'policy_loss': 28.0964, 'mean_log_prob': -0.54, 'mean_advantage': 51.8201, 'total_episodes': 528.0, 'total_timesteps': 42566.0}
[INFO] Iter 180, Step 45525, episodic return is 197.00. {'iteration': 180.0, 'performance': 201.0, 'ep_len': 201.0, 'ep_ret': 201.0, 'episode_len': 201.0, 'policy_loss': 30.9208, 'mean_log_prob': -0.5381, 'mean_advantage': 57.2793, 'total_episodes': 546.0, 'total_timesteps': 45525.0}
[INFO] Iter 180, episodic return 197.000 is greater than reward threshold 195.0. Congratulation! Now we exit the training process.
```

Environment is closed.

In [64]: # Run this cell without modification

```
pg_trainer_with_na, pg_result_with_na = run(PGTrainer, dict(
    learning_rate=0.001,
    max_episode_length=200,
    train_batch_size=200,
    env_name="CartPole-v1",
    normalize_advantage=True, # <<== Here!

    evaluate_interval=10,
    evaluate_num_episodes=10,
), 195.0)
```

```
[INFO] Iter 0, Step 207, episodic return is 20.50. {'iteration': 0.0, 'performance': 25.875, 'ep_len': 25.875, 'ep_ret': 25.875, 'episode_len': 207.0, 'policy_loss': 0.0028, 'mean_log_prob': -0.697, 'mean_advantage': 0.0, 'total_episodes': 8.0, 'total_timesteps': 207.0}
[INFO] Iter 10, Step 2362, episodic return is 33.70. {'iteration': 10.0, 'performance': 31.5714, 'ep_len': 31.5714, 'ep_ret': 31.5714, 'episode_len': 221.0, 'policy_loss': -0.013, 'mean_log_prob': -0.6819, 'mean_advantage': 0.0, 'total_episodes': 102.0, 'total_timesteps': 2362.0}
[INFO] Iter 20, Step 4631, episodic return is 40.60. {'iteration': 20.0, 'performance': 26.25, 'ep_len': 26.25, 'ep_ret': 26.25, 'episode_len': 210.0, 'policy_loss': -0.0144, 'mean_log_prob': -0.6628, 'mean_advantage': -0.0, 'total_episodes': 168.0, 'total_timesteps': 4631.0}
[INFO] Iter 30, Step 6856, episodic return is 52.30. {'iteration': 30.0, 'performance': 41.8, 'ep_len': 41.8, 'ep_ret': 41.8, 'episode_len': 209.0, 'policy_loss': -0.0055, 'mean_log_prob': -0.6289, 'mean_advantage': 0.0, 'total_episodes': 217.0, 'total_timesteps': 6856.0}
[INFO] Iter 40, Step 9218, episodic return is 76.70. {'iteration': 40.0, 'performance': 97.0, 'ep_len': 97.0, 'ep_ret': 97.0, 'episode_len': 291.0, 'policy_loss': -0.0009, 'mean_log_prob': -0.6125, 'mean_advantage': -0.0, 'total_episodes': 256.0, 'total_timesteps': 9218.0}
[INFO] Iter 50, Step 11952, episodic return is 76.90. {'iteration': 50.0, 'performance': 82.3333, 'ep_len': 82.3333, 'ep_ret': 82.3333, 'episode_len': 247.0, 'policy_loss': 0.0046, 'mean_log_prob': -0.5755, 'mean_advantage': 0.0, 'total_episodes': 289.0, 'total_timesteps': 11952.0}
[INFO] Iter 60, Step 14619, episodic return is 140.20. {'iteration': 60.0, 'performance': 135.0, 'ep_len': 135.0, 'ep_ret': 135.0, 'episode_len': 270.0, 'policy_loss': -0.0063, 'mean_log_prob': -0.5876, 'mean_advantage': -0.0, 'total_episodes': 317.0, 'total_timesteps': 14619.0}
[INFO] Iter 70, Step 17685, episodic return is 151.90. {'iteration': 70.0, 'performance': 201.0, 'ep_len': 201.0, 'ep_ret': 201.0, 'episode_len': 201.0, 'policy_loss': 0.0302, 'mean_log_prob': -0.5694, 'mean_advantage': -0.0, 'total_episodes': 339.0, 'total_timesteps': 17685.0}
[INFO] Iter 80, Step 20229, episodic return is 179.40. {'iteration': 80.0, 'performance': 201.0, 'ep_len': 201.0, 'ep_ret': 201.0, 'episode_len': 201.0, 'policy_loss': 0.0015, 'mean_log_prob': -0.5705, 'mean_advantage': -0.0, 'total_episodes': 355.0, 'total_timesteps': 20229.0}
[INFO] Iter 90, Step 23104, episodic return is 179.60. {'iteration': 90.0, 'performance': 201.0, 'ep_len': 201.0, 'ep_ret': 201.0, 'episode_len': 201.0, 'policy_loss': -0.0154, 'mean_log_prob': -0.5949, 'mean_advantage': -0.0, 'total_episodes': 374.0, 'total_timesteps': 23104.0}
[INFO] Iter 100, Step 25706, episodic return is 189.50. {'iteration': 100.0, 'performance': 174.0, 'ep_len': 174.0, 'ep_ret': 174.0, 'episode_len': 348.0, 'policy_loss': -0.0332, 'mean_log_prob': -0.5589, 'mean_advantage': -0.0, 'total_episodes': 390.0, 'total_timesteps': 25706.0}
[INFO] Iter 110, Step 28070, episodic return is 181.80. {'iteration': 110.0, 'performance': 201.0, 'ep_len': 201.0, 'ep_ret': 201.0, 'episode_len': 201.0, 'policy_loss': 0.0116, 'mean_log_prob': -0.5364, 'mean_advantage': -0.0, 'total_episodes': 403.0, 'total_timesteps': 28070.0}
[INFO] Iter 120, Step 30334, episodic return is 176.70. {'iteration': 120.0, 'performance': 113.0, 'ep_len': 113.0, 'ep_ret': 113.0, 'episode_len': 226.0, 'policy_loss': -0.0541, 'mean_log_prob': -0.5712, 'mean_advantage': 0.0, 'total_episodes': 416.0, 'total_timesteps': 30334.0}
[INFO] Iter 130, Step 32917, episodic return is 170.80. {'iteration': 130.0, 'performance': 159.0, 'ep_len': 159.0, 'ep_ret': 159.0, 'episode_len': 318.0, 'policy_loss': -0.0051, 'mean_log_prob': -0.5625, 'mean_advantage': -0.0, 'total_episodes': 430.0, 'total_timesteps': 32917.0}
[INFO] Iter 140, Step 35311, episodic return is 198.60. {'iteration': 140.0, 'performance': 201.0, 'ep_len': 201.0, 'ep_ret': 201.0, 'episode_len': 201.0, 'policy_loss': 0.0163, 'mean_log_prob': -0.5542, 'mean_advantage': -0.0, 'total_episodes': 443.0, 'total_timesteps': 35311.0}
```

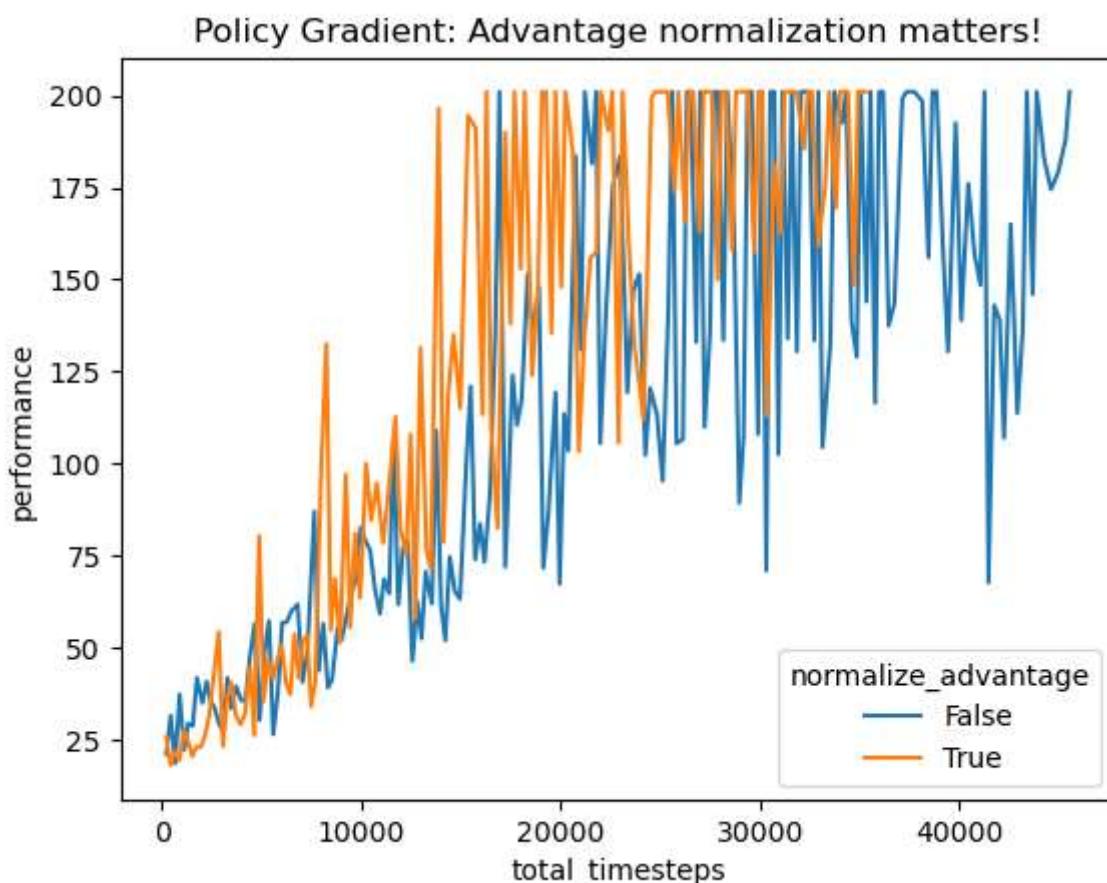
[INFO] Iter 140, episodic return 198.600 is greater than reward threshold 195.0. Congratulation! Now we exit the training process.
Environment is closed.

In [65]: # Run this cell without modification

```
pg_result_no_na_df = pd.DataFrame(pg_result_no_na)
pg_result_with_na_df = pd.DataFrame(pg_result_with_na)
pg_result_no_na_df["normalize_advantage"] = False
pg_result_with_na_df["normalize_advantage"] = True

ax = sns.lineplot(
    x="total_timesteps",
    y="performance",
    data=pd.concat([pg_result_no_na_df, pg_result_with_na_df]).reset_index(),
    hue="normalize_advantage"
)
ax.set_title("Policy Gradient: Advantage normalization matters!")
```

Out[65]: Text(0.5, 1.0, 'Policy Gradient: Advantage normalization matters!')



Section 4.4: Train REINFORCE in MetaDrive-Easy

In [66]: # Run this cell without modification

```
env_name = "MetaDrive-Tut-Easy-v0"

pg_trainer_metadrive_easy, pg_trainer_metadrive_easy_result = run(PGTrainer, dict(
    train_batch_size=2000,
    normalize_advantage=True,
    max_episode_length=200,
```

```

max_iteration=5000,
evaluate_interval=10,
evaluate_num_episodes=10,
learning_rate=0.001,
clip_norm=10.0,
env_name=env_name
), reward_threshold=120)

pg_trainer_metadrive_easy.save("pg_trainer_metadrive_easy.pt")

```

[INFO] Iter 0, Step 2010, episodic return is 2.50. {'iteration': 0.0, 'performance': 1.9263, 'ep_len': 201.0, 'ep_ret': 1.9263, 'episode_len': 2010.0, 'success_rate': 0.0, 'policy_loss': 0.0073, 'mean_log_prob': -2.1849, 'mean_advantage': 0.0, 'total_episodes': 10.0, 'total_timesteps': 2010.0}

[INFO] Iter 10, Step 22293, episodic return is 6.26. {'iteration': 10.0, 'performance': 7.8242, 'ep_len': 201.0, 'ep_ret': 7.8242, 'episode_len': 2010.0, 'success_rate': 0.0, 'policy_loss': -0.0409, 'mean_log_prob': -2.0147, 'mean_advantage': -0.0, 'total_episodes': 112.0, 'total_timesteps': 22293.0}

[INFO] Iter 20, Step 42863, episodic return is 67.18. {'iteration': 20.0, 'performance': 61.0495, 'ep_len': 103.95, 'ep_ret': 61.0495, 'episode_len': 2079.0, 'success_rate': 0.1, 'policy_loss': -0.0645, 'mean_log_prob': -1.1928, 'mean_advantage': 0.0, 'total_episodes': 264.0, 'total_timesteps': 42863.0}

[INFO] Iter 30, Step 63449, episodic return is 95.09. {'iteration': 30.0, 'performance': 97.9564, 'ep_len': 82.8, 'ep_ret': 97.9564, 'episode_len': 2070.0, 'success_rate': 0.6, 'policy_loss': -0.0366, 'mean_log_prob': -0.1513, 'mean_advantage': -0.0, 'total_episodes': 515.0, 'total_timesteps': 63449.0}

[INFO] Iter 40, Step 83812, episodic return is 122.04. {'iteration': 40.0, 'performance': 125.4393, 'ep_len': 92.0909, 'ep_ret': 125.4393, 'episode_len': 2026.0, 'success_rate': 1.0, 'policy_loss': 0.0033, 'mean_log_prob': -0.0171, 'mean_advantage': 0.0, 'total_episodes': 742.0, 'total_timesteps': 83812.0}

[INFO] Iter 40, episodic return 122.037 is greater than reward threshold 120. Congratulation! Now we exit the training process.

Environment is closed.

In [67]: # Run this cell without modification

```

# Render the Learned behavior
# NOTE: The Learned agent is marked by green color.
eval_reward, eval_info = evaluate(
    policy=pg_trainer_metadrive_easy.policy,
    num_episodes=1,
    env_name=pg_trainer_metadrive_easy.env_name,
    render="topdown", # Visualize the behaviors in top-down view
    verbose=True
)

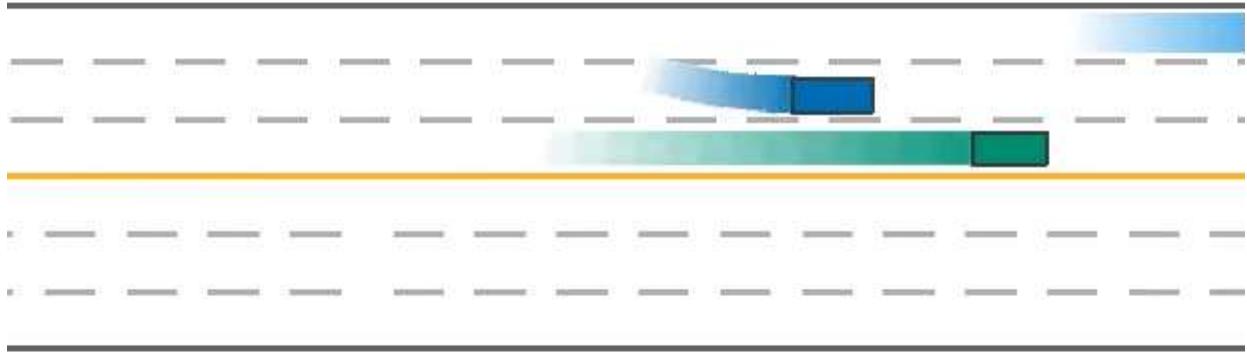
frames = [pygame.surfarray.array3d(f).swapaxes(0, 1) for f in eval_info["frames"]]

animate(frames)

print("REINFORCE agent achieves {} return in MetaDrive easy environment.".format(eval_

```

Evaluating 1/1 episodes. We are in 1/1000 steps. Current episode reward: 0.000
Evaluating 1/1 episodes. We are in 51/1000 steps. Current episode reward: 35.980



Section 5: Policy gradient with baseline

(20 / 100 points)

We compute the gradient of $Q = \mathbb{E} \sum_t r(a_t, s_t)$ w.r.t. the parameter to update the policy. Let's consider this case: when you take a so-so action that lead to positive expected return, the policy gradient is also positive and you will update your network toward this action. At the same time a potential better action is ignored.

To tackle this problem, we introduce the "baseline" when computing the policy gradient. The insight behind this is that we want to optimize the policy toward an action that are better than the "average action".

We introduce $b_t = \mathbb{E}_{a_t} \sum_{t'} \gamma^{t'-t} r(s_{t'}, a_{t'})$ as the baseline. It average the expected discount return of all possible actions at state s_t . So that the "advantage" achieved by action a_t can be evaluated via $\sum_{t'=t} \gamma^{t'-t} r(a_{t'}, s_{t'}) - b_t$

Therefore, the policy gradient becomes:

$$\nabla_\theta Q = \frac{1}{N} \sum_{i=1}^N \left[\left(\sum_t \nabla_\theta \log \pi_\theta(a_{i,t} | s_{i,t}) \left(\sum_{t'} \gamma^{t'-t} r(s_{i,t'}, a_{i,t'}) - b_{i,t} \right) \right) \right]$$

In our implementation, we estimate the baseline via an extra network `self.baseline`, which has same structure of policy network but output only a scalar value. We use the output of this network to serve as the baseline, while this network is updated by fitting the true value of expected return of current state: $\mathbb{E}_{a_t} \sum_{t'} \gamma^{t'-t} r(s_{t'}, a_{t'})$

The state-action values might have large variance if the reward function has large variance. It is not easy for a neural network to predict targets with large variance and extreme values. In implementation, we use a trick to match the distribution of baseline and values. During training, we first collect a batch of target values: $\{t_i = \mathbb{E}_{a_t} \sum_{t'} \gamma^{t'-t} r(s_{t'}, a_{t'})\}_i$. Then we normalize all targets to a standard distribution with mean = 0 and std = 1. Then we ask the baseline network to fit such normalized targets.

When computing the advantages, instead of using the output of baseline network as the baseline b , we firstly match the baseline distribution with state-action values, that is we "de-standardize" the baselines. The transformed baselines $b' = f(b)$ should has the same mean and STD with the action values.

After that, we compute the advantage of current action: $adv_{i,t} = \sum_{t'} \gamma^{t'-t} r(s_{i,t'}, a_{i,t'}) - b'_{i,t}$

By doing this, we mitigate the instability of training baseline.

Hint: We suggest to normalize an array via: `(x - x.mean()) / max(x.std(), 1e-6)`. The max term can mitigate numerical instability.

Section 5.1: Build PG method with baseline

```
In [74]: class PolicyGradientWithBaselineTrainer(PGTrainer):
    def initialize_parameters(self):
        # Build the actor in name of self.policy
        super().initialize_parameters()

        # TODO: Build the baseline network using PytorchModel class.
        self.baseline = PytorchModel(
            self.obs_dim,
            1,
            hidden_units=self.config["hidden_units"]
        )

        self.baseline_loss = nn.MSELoss()

        self.baseline_optimizer = torch.optim.Adam(
            self.baseline.parameters(),
            lr=self.config["learning_rate"]
        )

    def process_samples(self, samples):
        # Call the original process_samples function to get advantages
        tmp_samples, _ = super().process_samples(samples)
        values = tmp_samples["advantages"]
        samples["values"] = values # We add q_values into samples

        # Flatten the observations in all trajectories (still a numpy array)
        obs = np.concatenate(samples["obs"])

        assert obs.ndim == 2
        assert obs.shape[1] == self.obs_dim

        obs = self.to_tensor(obs)
        samples["flat_obs"] = obs

        # TODO: Compute the baseline by feeding observation to baseline network
        # Hint: baselines turns out to be a numpy array with the same shape of `values
        # that is: (batch size, )
        with torch.no_grad():
            baselines = self.baseline(obs).squeeze().cpu().numpy()

        assert baselines.shape == values.shape

        # TODO: Match the distribution of baselines to the values.
        # Hint: We expect to see baselines.std almost equals to values.std,
        # and baselines.mean almost equals to values.mean.
        if self.config["normalize_advantage"]:
            # Normalize the baseline to match the distribution of values
            mean_values = np.mean(values)
            std_values = np.std(values)

            mean_baselines = np.mean(baselines)
            std_baselines = np.std(baselines)
```

```

        baselines = (baselines - mean_baselines) * (std_values / std_baselines) +

    # Compute the advantage
    advantages = values - baselines
    samples["advantages"] = advantages
    process_info = {"mean_baseline": float(np.mean(baselines))}
    return samples, process_info

def update_network(self, processed_samples):
    update_info = super().update_network(processed_samples)
    update_info.update(self.update_baseline(processed_samples))
    return update_info

def update_baseline(self, processed_samples):
    self.baseline.train()
    obs = processed_samples["flat_obs"]

    # TODO: Normalize `values` to have mean=0, std=1.
    values = processed_samples["values"]
    mean_values = np.mean(values)
    std_values = np.std(values)
    values = (values - mean_values) / std_values

    values = self.to_tensor(values[:, np.newaxis])

    baselines = self.baseline(obs)

    self.baseline_optimizer.zero_grad()
    loss = self.baseline_loss(input=baselines, target=values)
    loss.backward()

    # Clip the gradient
    torch.nn.utils.clip_grad_norm_(
        self.baseline.parameters(), self.config["clip_gradient"]
    )

    self.baseline_optimizer.step()
    self.baseline.eval()
    return dict(baseline_loss=loss.item())

```

Section 5.2: Run PG w/ baseline in CartPole

In [75]: # Run this cell without modification

```

pg_trainer_wb_cartpole, pg_trainer_wb_cartpole_result = run(PolicyGradientWithBaseline
    learning_rate=0.001,
    max_episode_length=200,
    train_batch_size=200,

    env_name="CartPole-v1",
    normalize_advantage=True,

    evaluate_interval=10,
    evaluate_num_episodes=10,
), 195.0

```

```
[INFO] Iter 0, Step 213, episodic return is 20.20. {'iteration': 0.0, 'performance': 19.3636, 'ep_len': 19.3636, 'ep_ret': 19.3636, 'episode_len': 213.0, 'mean_baseline': -0.0, 'policy_loss': -0.0001, 'mean_log_prob': -0.6947, 'mean_advantage': 0.0, 'baseline_loss': 1.0078, 'total_episodes': 11.0, 'total_timesteps': 213.0}
[INFO] Iter 10, Step 2357, episodic return is 31.20. {'iteration': 10.0, 'performance': 25.4444, 'ep_len': 25.4444, 'ep_ret': 25.4444, 'episode_len': 229.0, 'mean_baseline': 0.0, 'policy_loss': -0.0296, 'mean_log_prob': -0.6859, 'mean_advantage': 0.0, 'baseline_loss': 0.9589, 'total_episodes': 106.0, 'total_timesteps': 2357.0}
[INFO] Iter 20, Step 4546, episodic return is 34.10. {'iteration': 20.0, 'performance': 55.25, 'ep_len': 55.25, 'ep_ret': 55.25, 'episode_len': 221.0, 'mean_baseline': -0.0, 'policy_loss': -0.0319, 'mean_log_prob': -0.6689, 'mean_advantage': -0.0, 'baseline_loss': 0.929, 'total_episodes': 177.0, 'total_timesteps': 4546.0}
[INFO] Iter 30, Step 6769, episodic return is 45.50. {'iteration': 30.0, 'performance': 44.8, 'ep_len': 44.8, 'ep_ret': 44.8, 'episode_len': 224.0, 'mean_baseline': -0.0, 'policy_loss': -0.0767, 'mean_log_prob': -0.635, 'mean_advantage': -0.0, 'baseline_loss': 0.8551, 'total_episodes': 230.0, 'total_timesteps': 6769.0}
[INFO] Iter 40, Step 9074, episodic return is 65.80. {'iteration': 40.0, 'performance': 68.6667, 'ep_len': 68.6667, 'ep_ret': 68.6667, 'episode_len': 206.0, 'mean_baseline': 0.0, 'policy_loss': -0.0598, 'mean_log_prob': -0.6212, 'mean_advantage': 0.0, 'baseline_loss': 0.8951, 'total_episodes': 274.0, 'total_timesteps': 9074.0}
[INFO] Iter 50, Step 11377, episodic return is 64.20. {'iteration': 50.0, 'performance': 50.8, 'ep_len': 50.8, 'ep_ret': 50.8, 'episode_len': 254.0, 'mean_baseline': -0.0, 'policy_loss': -0.0839, 'mean_log_prob': -0.5929, 'mean_advantage': 0.0, 'baseline_loss': 0.8142, 'total_episodes': 311.0, 'total_timesteps': 11377.0}
[INFO] Iter 60, Step 13600, episodic return is 79.70. {'iteration': 60.0, 'performance': 67.0, 'ep_len': 67.0, 'ep_ret': 67.0, 'episode_len': 201.0, 'mean_baseline': 0.0, 'policy_loss': -0.0452, 'mean_log_prob': -0.5666, 'mean_advantage': 0.0, 'baseline_loss': 0.6276, 'total_episodes': 346.0, 'total_timesteps': 13600.0}
[INFO] Iter 70, Step 16174, episodic return is 78.10. {'iteration': 70.0, 'performance': 96.0, 'ep_len': 96.0, 'ep_ret': 96.0, 'episode_len': 288.0, 'mean_baseline': 0.0, 'policy_loss': -0.0353, 'mean_log_prob': -0.588, 'mean_advantage': -0.0, 'baseline_loss': 0.4304, 'total_episodes': 379.0, 'total_timesteps': 16174.0}
[INFO] Iter 80, Step 18481, episodic return is 61.30. {'iteration': 80.0, 'performance': 91.0, 'ep_len': 91.0, 'ep_ret': 91.0, 'episode_len': 273.0, 'mean_baseline': 0.0, 'policy_loss': -0.0139, 'mean_log_prob': -0.5768, 'mean_advantage': -0.0, 'baseline_loss': 0.4168, 'total_episodes': 409.0, 'total_timesteps': 18481.0}
[INFO] Iter 90, Step 20800, episodic return is 62.50. {'iteration': 90.0, 'performance': 47.8, 'ep_len': 47.8, 'ep_ret': 47.8, 'episode_len': 239.0, 'mean_baseline': 0.0, 'policy_loss': -0.0388, 'mean_log_prob': -0.5763, 'mean_advantage': 0.0, 'baseline_loss': 0.4871, 'total_episodes': 443.0, 'total_timesteps': 20800.0}
[INFO] Iter 100, Step 23264, episodic return is 98.70. {'iteration': 100.0, 'performance': 105.0, 'ep_len': 105.0, 'ep_ret': 105.0, 'episode_len': 210.0, 'mean_baseline': -0.0, 'policy_loss': -0.0023, 'mean_log_prob': -0.5229, 'mean_advantage': 0.0, 'baseline_loss': 0.2842, 'total_episodes': 468.0, 'total_timesteps': 23264.0}
[INFO] Iter 110, Step 25896, episodic return is 124.90. {'iteration': 110.0, 'performance': 81.6667, 'ep_len': 81.6667, 'ep_ret': 81.6667, 'episode_len': 245.0, 'mean_baseline': -0.0, 'policy_loss': -0.0091, 'mean_log_prob': -0.5462, 'mean_advantage': 0.0, 'baseline_loss': 0.1446, 'total_episodes': 493.0, 'total_timesteps': 25896.0}
[INFO] Iter 120, Step 28695, episodic return is 159.00. {'iteration': 120.0, 'performance': 110.0, 'ep_len': 110.0, 'ep_ret': 110.0, 'episode_len': 330.0, 'mean_baseline': 0.0, 'policy_loss': -0.0324, 'mean_log_prob': -0.5467, 'mean_advantage': -0.0, 'baseline_loss': 0.5349, 'total_episodes': 517.0, 'total_timesteps': 28695.0}
[INFO] Iter 130, Step 31380, episodic return is 166.80. {'iteration': 130.0, 'performance': 107.0, 'ep_len': 107.0, 'ep_ret': 107.0, 'episode_len': 321.0, 'mean_baseline': -0.0, 'policy_loss': -0.007, 'mean_log_prob': -0.5486, 'mean_advantage': 0.0, 'baseline_loss': 0.2059, 'total_episodes': 540.0, 'total_timesteps': 31380.0}
[INFO] Iter 140, Step 34098, episodic return is 152.70. {'iteration': 140.0, 'performance': 153.0, 'ep_len': 153.0, 'ep_ret': 153.0, 'episode_len': 306.0, 'mean_baseline': 0.0, 'policy_loss': -0.0028, 'mean_log_prob': -0.526, 'mean_advantage': -0.0, 'baseline_loss': 0.619, 'total_episodes': 556.0, 'total_timesteps': 34098.0}
```

```
[INFO] Iter 150, Step 36935, episodic return is 166.90. {'iteration': 150.0, 'performance': 151.0, 'ep_len': 151.0, 'ep_ret': 151.0, 'episode_len': 302.0, 'mean_baseline': -0.0, 'policy_loss': -0.0052, 'mean_log_prob': -0.5098, 'mean_advantage': -0.0, 'baseline_loss': 0.3269, 'total_episodes': 578.0, 'total_timesteps': 36935.0}
[INFO] Iter 160, Step 39235, episodic return is 199.10. {'iteration': 160.0, 'performance': 201.0, 'ep_len': 201.0, 'ep_ret': 201.0, 'episode_len': 201.0, 'mean_baseline': 0.0, 'policy_loss': -0.0102, 'mean_log_prob': -0.5065, 'mean_advantage': -0.0, 'baseline_loss': 0.2892, 'total_episodes': 591.0, 'total_timesteps': 39235.0}
[INFO] Iter 160, episodic return 199.100 is greater than reward threshold 195.0. Congratulation! Now we exit the training process.
```

Environment is closed.

Section 5.3: Run PG w/ baseline in MetaDrive-Easy

In [76]: # Run this cell without modification

```
env_name = "MetaDrive-Tut-Easy-v0"

pg_trainer_wb_metadrive_easy, pg_trainer_wb_metadrive_easy_result = run(
    PolicyGradientWithBaselineTrainer,
    dict(
        train_batch_size=2000,
        normalize_advantage=True,
        max_episode_length=200,
        max_iteration=5000,
        evaluate_interval=10,
        evaluate_num_episodes=10,
        learning_rate=0.001,
        clip_norm=10.0,
        env_name=env_name
    ),
    reward_threshold=120
)

pg_trainer_wb_metadrive_easy.save("pg_trainer_wb_metadrive_easy.pt")
```

```
[INFO] Iter 0, Step 2010, episodic return is 3.27. {'iteration': 0.0, 'performance': 2.7025, 'ep_len': 201.0, 'ep_ret': 2.7025, 'episode_len': 2010.0, 'success_rate': 0.0, 'mean_baseline': 0.0, 'policy_loss': -0.0018, 'mean_log_prob': -2.1923, 'mean_advantage': -0.0, 'baseline_loss': 1.0498, 'total_episodes': 10.0, 'total_timesteps': 2010.0}
[INFO] Iter 10, Step 22429, episodic return is 8.21. {'iteration': 10.0, 'performance': 9.2284, 'ep_len': 201.0, 'ep_ret': 9.2284, 'episode_len': 2010.0, 'success_rate': 0.0, 'mean_baseline': 0.0, 'policy_loss': -0.0383, 'mean_log_prob': -2.0032, 'mean_advantage': 0.0, 'baseline_loss': 0.998, 'total_episodes': 112.0, 'total_timesteps': 22429.0}
[INFO] Iter 20, Step 43369, episodic return is 38.00. {'iteration': 20.0, 'performance': 49.9573, 'ep_len': 122.1765, 'ep_ret': 49.9573, 'episode_len': 2077.0, 'success_rate': 0.1765, 'mean_baseline': 0.0, 'policy_loss': -0.0599, 'mean_log_prob': -1.4524, 'mean_advantage': -0.0, 'baseline_loss': 0.9985, 'total_episodes': 250.0, 'total_time_steps': 43369.0}
[INFO] Iter 30, Step 63722, episodic return is 82.84. {'iteration': 30.0, 'performance': 80.4937, 'ep_len': 77.8077, 'ep_ret': 80.4937, 'episode_len': 2023.0, 'success_rate': 0.3846, 'mean_baseline': 0.0, 'policy_loss': 0.0225, 'mean_log_prob': -0.262, 'mean_advantage': -0.0, 'baseline_loss': 0.9859, 'total_episodes': 508.0, 'total_time_steps': 63722.0}
[INFO] Iter 40, Step 84134, episodic return is 125.61. {'iteration': 40.0, 'performance': 123.4515, 'ep_len': 91.4091, 'ep_ret': 123.4515, 'episode_len': 2011.0, 'success_rate': 0.9545, 'mean_baseline': -0.0, 'policy_loss': -0.005, 'mean_log_prob': -0.0058, 'mean_advantage': 0.0, 'baseline_loss': 0.8536, 'total_episodes': 743.0, 'total_time_steps': 84134.0}
[INFO] Iter 40, episodic return 125.608 is greater than reward threshold 120. Congratulation! Now we exit the training process.
```

Environment is closed.

```
In [77]: # Run this cell without modification

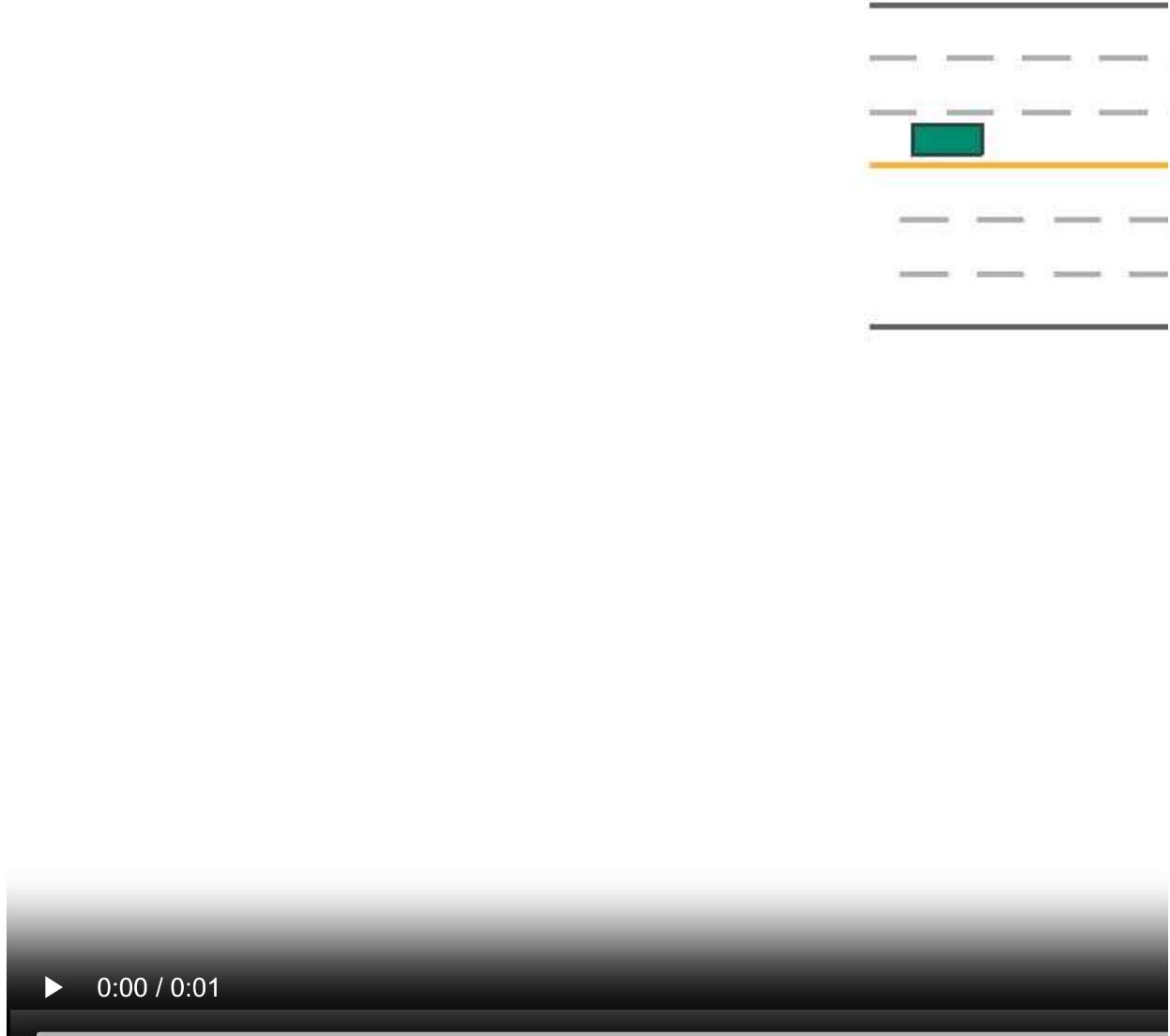
# Render the Learned behavior
# NOTE: The learned agent is marked by green color.
eval_reward, eval_info = evaluate(
    policy=pg_trainer_wb_metadrive_easy.policy,
    num_episodes=1,
    env_name=pg_trainer_wb_metadrive_easy.env_name,
    render="topdown", # Visualize the behaviors in top-down view
    verbose=True
)

frames = [pygame.surfarray.array3d(f).swapaxes(0, 1) for f in eval_info["frames"]]

print(
    "PG agent achieves {} return and {} success rate in MetaDrive easy environment.".format(
        eval_reward, eval_info["success_rate"]
    )
)

animate(frames)
```

Evaluating 1/1 episodes. We are in 1/1000 steps. Current episode reward: 0.000
Evaluating 1/1 episodes. We are in 51/1000 steps. Current episode reward: 35.980
PG agent achieves 125.53851204681443 return and 1.0 success rate in MetaDrive easy environment.



Section 5.4: Run PG with baseline in MetaDrive-Hard

The minimum goal to is to achieve episodic return > 20, which costs nearly 20 iterations and ~100k steps.

Bonus

BONUS can be earned if you can improve the training performance by adjusting hyper-parameters and optimizing code. Improvement means achieving > 0.0 success rate.

However, I can't guarantee it is feasible to solve this task with PG via simply tweaking the hyper-parameters more carefully. Please creates a independent markdown cell to highlight your improvement.

```
In [79]: # Run this cell without modification

env_name = "MetaDrive-Tut-Hard-v0"

pg_trainer_wb_metadrive_hard, pg_trainer_wb_metadrive_hard_result = run(
    PolicyGradientWithBaselineTrainer,
    dict(
        train_batch_size=4000,
        normalize_advantage=True,
        max_episode_length=1000,
        max_iteration=5000,
        evaluate_interval=5,
        evaluate_num_episodes=10,
        learning_rate=0.001,
        clip_norm=10.0,
        env_name=env_name
    ),
    reward_threshold=40 # We just set the reward threshold to 20. Feel free to adjust
)
pg_trainer_wb_metadrive_hard.save("pg_trainer_wb_metadrive_hard.pt")
```

```
[INFO] Iter 0, Step 4898, episodic return is 11.89. {'iteration': 0.0, 'performance': 12.4349, 'ep_len': 979.6, 'ep_ret': 12.4349, 'episode_len': 4898.0, 'success_rate': 0.0, 'mean_baseline': 0.0, 'policy_loss': -0.001, 'mean_log_prob': -3.2157, 'mean_advantage': -0.0, 'baseline_loss': 1.0067, 'total_episodes': 5.0, 'total_timesteps': 4898.0}
[INFO] Iter 5, Step 27731, episodic return is 13.87. {'iteration': 5.0, 'performance': 14.7524, 'ep_len': 966.2, 'ep_ret': 14.7524, 'episode_len': 4831.0, 'success_rate': 0.0, 'mean_baseline': 0.0, 'policy_loss': -0.0064, 'mean_log_prob': -3.1976, 'mean_advantage': -0.0, 'baseline_loss': 0.9989, 'total_episodes': 30.0, 'total_timesteps': 27731.0}
[INFO] Iter 10, Step 49553, episodic return is 16.50. {'iteration': 10.0, 'performance': 22.5026, 'ep_len': 1001.0, 'ep_ret': 22.5026, 'episode_len': 4004.0, 'success_rate': 0.0, 'mean_baseline': 0.0, 'policy_loss': -0.0204, 'mean_log_prob': -3.1607, 'mean_advantage': -0.0, 'baseline_loss': 0.9998, 'total_episodes': 54.0, 'total_timesteps': 49553.0}
[INFO] Iter 15, Step 71853, episodic return is 22.12. {'iteration': 15.0, 'performance': 13.8197, 'ep_len': 516.375, 'ep_ret': 13.8197, 'episode_len': 4131.0, 'success_rate': 0.0, 'mean_baseline': -0.0, 'policy_loss': -0.017, 'mean_log_prob': -3.0339, 'mean_advantage': 0.0, 'baseline_loss': 1.0022, 'total_episodes': 88.0, 'total_timesteps': 71853.0}
[INFO] Iter 20, Step 93907, episodic return is 36.86. {'iteration': 20.0, 'performance': 54.7149, 'ep_len': 568.25, 'ep_ret': 54.7149, 'episode_len': 4546.0, 'success_rate': 0.0, 'mean_baseline': -0.0, 'policy_loss': -0.0585, 'mean_log_prob': -2.6196, 'mean_advantage': 0.0, 'baseline_loss': 1.0022, 'total_episodes': 128.0, 'total_timesteps': 93907.0}
[INFO] Iter 25, Step 114266, episodic return is 43.14. {'iteration': 25.0, 'performance': 40.8311, 'ep_len': 72.8364, 'ep_ret': 40.8311, 'episode_len': 4006.0, 'success_rate': 0.0182, 'mean_baseline': -0.0, 'policy_loss': -0.0078, 'mean_log_prob': -1.4459, 'mean_advantage': -0.0, 'baseline_loss': 0.9952, 'total_episodes': 323.0, 'total_timesteps': 114266.0}
[INFO] Iter 25, episodic return 43.145 is greater than reward threshold 40. Congratulations! Now we exit the training process.
```

Environment is closed.

Modified the reward_threshold to be 40 instead of 20. In metadrive_easy, it reached a success rate greater than 0 at approximately episodic return of 38.

```
In [80]: # Run this cell without modification

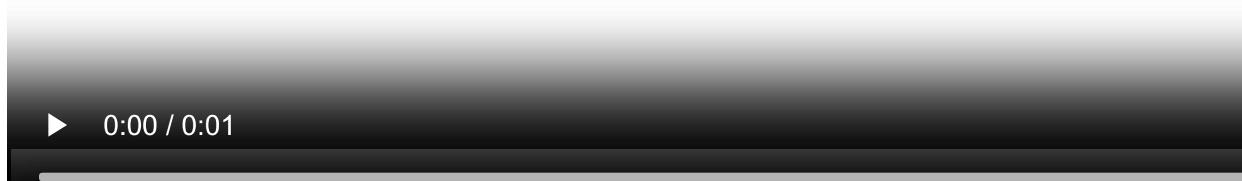
# Render the Learned behavior
# NOTE: The Learned agent is marked by green color.
eval_reward, eval_info = evaluate(
    policy=pg_trainer_wb_metadrive_hard.policy,
    num_episodes=10,
    env_name=pg_trainer_wb_metadrive_hard.env_name,
    render=None,
    verbose=False
)

_, eval_info_render = evaluate(
    policy=pg_trainer_wb_metadrive_hard.policy,
    num_episodes=1,
    env_name=pg_trainer_wb_metadrive_hard.env_name,
    render="topdown", # Visualize the behaviors in top-down view
    verbose=True
)

frames = [pygame.surfarray.array3d(f).swapaxes(0, 1) for f in eval_info_render["frames"]]
```

```
print(  
    "PG agent achieves {} return and {} success rate in MetaDrive easy environment.".f  
        eval_reward, eval_info["success_rate"]  
)  
  
animate(frames)
```

Evaluating 1/1 episodes. We are in 1/1000 steps. Current episode reward: 0.000
Evaluating 1/1 episodes. We are in 51/1000 steps. Current episode reward: 25.997
PG agent achieves 54.974254142096264 return and 0.1 success rate in MetaDrive easy en
vironment.



Conclusion

In this assignment, we learn how to build naive Q learning, Deep Q Network and Policy Gradient methods.

Following the submission instruction in the assignment to submit your assignment. Thank you!
