**Install necessary Libraries**

```
1  !pip install gymnasium[mujoco]
```

```
Requirement already satisfied: gymnasium[mujoco] in /usr/local/lib/python3.10/dist-packages (0.29.1)
Requirement already satisfied: numpy>=1.21.0 in /usr/local/lib/python3.10/dist-packages (from gymnasium[mujoco]) (1.23.5)
Requirement already satisfied: cloudpickle>=1.2.0 in /usr/local/lib/python3.10/dist-packages (from gymnasium[mujoco]) (2.2.1)
Requirement already satisfied: typing-extensions>=4.3.0 in /usr/local/lib/python3.10/dist-packages (from gymnasium[mujoco]) (4.5.0)
Requirement already satisfied: farama-notifications>=0.0.1 in /usr/local/lib/python3.10/dist-packages (from gymnasium[mujoco]) (0.0.4)
Requirement already satisfied: mujoco>=2.3.3 in /usr/local/lib/python3.10/dist-packages (from gymnasium[mujoco]) (3.1.1)
Requirement already satisfied: imageio>=2.14.1 in /usr/local/lib/python3.10/dist-packages (from gymnasium[mujoco]) (2.31.6)
Requirement already satisfied: pillow<10.1.0,>=8.3.2 in /usr/local/lib/python3.10/dist-packages (from imageio>=2.14.1->gymnasium[mujoco]
Requirement already satisfied: absl-py in /usr/local/lib/python3.10/dist-packages (from mujoco>=2.3.3->gymnasium[mujoco]) (1.4.0)
Requirement already satisfied: etils[epath] in /usr/local/lib/python3.10/dist-packages (from mujoco>=2.3.3->gymnasium[mujoco]) (1.6.0)
Requirement already satisfied: glfw in /usr/local/lib/python3.10/dist-packages (from mujoco>=2.3.3->gymnasium[mujoco]) (2.6.4)
Requirement already satisfied: pyopengl in /usr/local/lib/python3.10/dist-packages (from mujoco>=2.3.3->gymnasium[mujoco]) (3.1.7)
Requirement already satisfied: fsspec in /usr/local/lib/python3.10/dist-packages (from etils[epath]->mujoco>=2.3.3->gymnasium[mujoco]) (
Requirement already satisfied: importlib_resources in /usr/local/lib/python3.10/dist-packages (from etils[epath]->mujoco>=2.3.3->gymnasi
Requirement already satisfied: zipp in /usr/local/lib/python3.10/dist-packages (from etils[epath]->mujoco>=2.3.3->gymnasium[mujoco]) (3.
```

**Import Necessary Libraries**

```
1   from __future__ import annotations
2
3   import sys
4   import random
5   import argparse
6   from collections import deque
7
8   import matplotlib.pyplot as plt
9   import numpy as np
10  import pandas as pd
11  import seaborn as sns
12  import torch
13  import torch.nn as nn
14  from torch.distributions.normal import Normal
15
16  import gymnasium as gym
17
```

**Define Policy Network for continuous action space**

```
1  class Policy_Network(nn.Module):
2      """Parametrized Policy Network."""
3
4      def __init__(self, obs_space_dims: int, action_space_dims: int):
5          """Creates a neural network responsible for predicting the mean and
6          standard deviation of a normal distribution from which actions are
7          subsequently sampled.
8
9          Args:
10             obs_space_dims: The dimensions of the observation space
11             action_space_dims: The dimensions of the action space
12         """
13         super().__init__()
14
15         hidden_space1 = 16
16         hidden_space2 = 32
17
18         # Shared Network
19         self.shared_net = nn.Sequential(
20             nn.Linear(obs_space_dims, hidden_space1),
21             nn.Tanh(),
22             nn.Linear(hidden_space1, hidden_space2),
23             nn.Tanh(),
24         )
25
26         # Policy Mean specific Linear Layer
27         self.policy_mean_net = nn.Sequential(
28             nn.Linear(hidden_space2, action_space_dims)
29         )
30
31         # Initialize log_std as a parameter
32         self.log_std = nn.Parameter(torch.zeros(action_space_dims))
33
34     def forward(self, x: torch.Tensor) -> tuple[torch.Tensor, torch.Tensor]:
35         """Given an observation, this function provides the mean and standard deviation for sampling an action from a normal distribution
36
37             Arguments:
38                 x: The observation obtained from the environment
39
40             Returns:
41                 action_means: The predicted mean of the normal distribution
42                 action_stddevs: The predicted standard deviation of the normal distribution
43         """
44         shared_features = self.shared_net(x.float())
45
46         action_means = self.policy_mean_net(shared_features)
47         action_stddevs = self.log_std.exp()
48
49         return action_means, action_stddevs
```

### Define Value Network

```
1  class ValueNetwork(nn.Module):
2      def __init__(self, input_size, hidden_size):
3          super(ValueNetwork, self).__init__()
4          self.fc = nn.Sequential(
5              nn.Linear(input_size, hidden_size),
6              nn.ReLU(),
7              nn.Linear(hidden_size, hidden_size),
8              nn.ReLU(),
9              nn.Linear(hidden_size, 1)
10         )
11
12     def forward(self, x):
13         return self.fc(x)
```

### Define PolicyGradientAgent

```python
1  class PolicyGradientAgent:
2      """Unified Policy Gradient algorithm that can work with or without a baseline."""
3
4      def __init__(self, env, obs_space_dims, action_space_dims, lr=1e-5, use_baseline=False):
5          """Initializes an agent that learns a policy via REINFORCE or PGB algorithm.
6
7          Args:
8              obs_space_dims: Dimension of the observation space.
9              action_space_dims: Dimension of the action space.
10             lr: Learning rate for policy optimization.
11             use_baseline: Whether to use a baseline (value network).
12         """
13         self.learning_rate = lr
14         self.gamma = 0.99
15         self.eps = 1e-6
16         self.use_baseline = use_baseline
17         self.env = env
18
19         self.probs = []
20         self.rewards = []
21         self.states = [] if use_baseline else None
22
23         self.policy_net = Policy_Network(obs_space_dims, action_space_dims)
24         self.policy_optimizer = torch.optim.AdamW(self.policy_net.parameters(), lr=self.learning_rate)
25
26         if use_baseline:
27             self.value_net = ValueNetwork(obs_space_dims, 16)
28             self.value_optimizer = torch.optim.Adam(self.value_net.parameters(), lr=1e-5)
29
30     def sample_action(self, state: np.ndarray) -> np.ndarray:
31         state_tensor = torch.tensor(np.array([state]), dtype=torch.float32)
32         action_means, action_stddevs = self.policy_net(state_tensor)
33         distrib = Normal(action_means[0], action_stddevs[0].clamp(min=self.eps))
34         action = distrib.sample()
35         prob = distrib.log_prob(action)
36         self.probs.append(prob)
37         return action.clamp(min=self.env.action_space.low[0], max=self.env.action_space.high[0]).numpy(
38
39     def update(self):
40         # Calculate reward-to-go for all time steps
41         R = 0
42         deltas = []
43         for r in self.rewards[::-1]:
44             R = r + self.gamma * R
45             deltas.insert(0, R)
46
47         # Convert list to tensor
48         deltas = torch.tensor(deltas, dtype=torch.float32)
49
50         # Calculate policy loss
51         policy_loss = -torch.stack(self.probs).mean() * deltas
52
53         # Update the policy network
54         self.policy_optimizer.zero_grad()
55         # Sum or average the elements to create a scalar loss
56         scalar_policy_loss = policy_loss.sum()  # or policy_loss.mean()
57
58         # Now you can call backward on the scalar loss
59         scalar_policy_loss.backward()
60         self.policy_optimizer.step()
61
62         # If using baseline, calculate value loss and update value network
63         if self.use_baseline:
64             states_tensor = torch.tensor(np.array(self.states), dtype=torch.float32)
65             values = self.value_net(states_tensor).squeeze()
66             advantages = deltas - values
67
68             value_loss = (advantages ** 2).mean()
69             self.value_optimizer.zero_grad()
70             # # Sum or average the elements to create a scalar loss
71             # scalar_value_loss = policy_loss.sum()  # or policy_loss.mean()
72
73             # # Now you can call backward on the scalar loss
74             # scalar_value_loss.backward()
75             value_loss.backward()
76             self.value_optimizer.step()
77
```

```
78        # Reset episode data
79        self.probs = []
80        self.rewards = []
81        if self.use_baseline:
82            self.states = []
```

## Training Loop-Function-Vanilla Policy Gradient

```python
1 # Define a function to run the Vanilla Policy Gradient (REINFORCE) algorithm
2 def run_vanilla_pg(env, total_num_episodes, seeds, obs_space_dims, action_space_dims):
3     rewards_over_seeds = []
4
5     for seed in seeds:
6         torch.manual_seed(seed)
7         random.seed(seed)
8         np.random.seed(seed)
9
10        agent = PolicyGradientAgent(env, obs_space_dims, action_space_dims, use_baseline=False)  # Initialize the REINFORCE agent
11        reward_over_episodes = []
12
13        for episode in range(total_num_episodes):
14            obs, info = env.reset()
15            done = False
16            while not done:
17                action = agent.sample_action(obs)
18                obs, reward, terminated, truncated, info = env.step(action)
19                agent.rewards.append(reward)
20                done = terminated or truncated
21
22            reward_over_episodes.append(sum(agent.rewards))
23            agent.update()
24
25            if episode % 1000 == 0:
26                avg_reward = np.mean(reward_over_episodes[-100:])  # Calculate average of last 100 episodes
27                print("Episode: ", episode, "Average Reward: ", avg_reward)
28
29        rewards_over_seeds.append(reward_over_episodes)
30
31    return rewards_over_seeds
```

## Training Loop-Function-Policy Gradient with Baseline

```
1 # Define a function to run the Policy Gradient with Baseline algorithm
2 def run_pgb(env, total_num_episodes, seeds, obs_space_dims, action_space_dims, learning_rate):
3     rewards_over_seeds = []
4
5     for seed in seeds:
6         torch.manual_seed(seed)
7         random.seed(seed)
8         np.random.seed(seed)
9
10        agent = PolicyGradientAgent(env, obs_space_dims, action_space_dims,
11                              lr=learning_rate, use_baseline=True)  # Initialize the PGB agent
12        reward_over_episodes = []
13
14        for episode in range(total_num_episodes):
15            obs, info = env.reset()
16            done = False
17            while not done:
18                action = agent.sample_action(obs)
19                agent.states.append(obs)
20                obs, reward, terminated, truncated, info = env.step(action)
21                agent.rewards.append(reward)
22                done = terminated or truncated
23
24            reward_over_episodes.append(sum(agent.rewards))
25            agent.update()
26
27            if episode % 1000 == 0:
28                avg_reward = np.mean(reward_over_episodes[-100:])  # Calculate average of last 100 episodes
29                print("Episode: ", episode, "Average Reward: ", avg_reward)
30
31        rewards_over_seeds.append(reward_over_episodes)
32
33    return rewards_over_seeds
```

**Plotting**

⌄  rolling average function

```
1 def rolling_average(data, window_size=10):
2     """Compute rolling average for smoother plots."""
3     return pd.Series(data).rolling(window=window_size).mean()
```

```
1 def plot_learning_curve(smoothed_rewards, label, title, xlabel, ylabel):
2     plt.figure(figsize=(12, 8))
3     plt.plot(smoothed_rewards, label=label)
4     plt.title(title)
5     plt.xlabel(xlabel)
6     plt.ylabel(ylabel)
7     plt.legend()
8     plt.grid(True)
9     plt.show()
```

```
1 def plot_combined_learning_curve(smoothed_rewards, label, title, xlabel, ylabel):
2     plt.figure(figsize=(12, 8))
3     plt.plot(smoothed_rewards, label=label)
4     plt.title(title)
5     plt.xlabel(xlabel)
6     plt.ylabel(ylabel)
7     plt.legend()
8     plt.grid(True)
9     plt.show()
```

## ⌄ Main function

```python
1 def main():
2     parser = argparse.ArgumentParser()
3     parser.add_argument('--algo', type=str, default='pg', choices=['pg', 'pgb'],
4                         help='Algorithm to run: "pg" for vanilla policy gradient, "pgb" for policy gradient with baseline')
5     args = parser.parse_args()
6
7     env = gym.make("Ant-v4")
8     total_num_episodes = 5000
9     seeds = [8]  # List of seeds you want to use
10    obs_space_dims = env.observation_space.shape[0]
11    action_space_dims = env.action_space.shape[0]
12
13    if args.algo == 'pg':
14        # Run Vanilla Policy Gradient
15        rewards_vpg = run_vanilla_pg(env, total_num_episodes, seeds, obs_space_dims, action_space_dims)
16
17        # Flatten the list to have a single list of rewards over all episodes
18        flattened_rewards = [reward for seed_rewards in rewards_vpg for reward in seed_rewards]
19
20        # Compute rolling average
21        smoothed_rewards = rolling_average(flattened_rewards, window_size=50)
22
23        plot_learning_curve( smoothed_rewards, 'Vanilla Policy Gradient',
24                            'Learning Curve for Vanilla Policy Gradient Algorithm on Ant-v4',
25                            'Episodes',
26                            'Undiscounted Return'
27                            )
28        ### save variable smoothed_rewards
29
30    elif args.algo == 'pgb':
31        # Run Policy Gradient with Baseline
32        rewards_pgb = run_pgb(env, total_num_episodes, seeds, obs_space_dims,
33                              action_space_dims, learning_rate= 1e-5)
34
35        # Flatten the list to have a single list of rewards over all episodes
36        flattened_rewards = [reward for seed_rewards in rewards_pgb for reward in seed_rewards]
37
38        # Compute rolling average
39        smoothed_rewards = rolling_average(flattened_rewards, window_size=50)
40
41        plot_learning_curve( smoothed_rewards, 'Policy Gradient with Baseline',
42                            'Learning Curve for Policy Gradient with Baseline Algorithm on Ant-v4',
43                            'Episodes',
44                            'Undiscounted Return'
45                            )
46
47        # Question 3
48        learning_rates = [1e-4, 1e-5, 1e-6]
49        results = {}
50
51        for lr in learning_rates:
52            print(f"Running experiment with learning rate: {lr}")
53            rewards_pgb = run_pgb(env, total_num_episodes, seeds, obs_space_dims,
54                                  action_space_dims, learning_rate= lr)
55            flattened_reward = [reward for seed_rewards in rewards_pgb for reward in seed_rewards]
56            results[f"lr={lr}"] = flattened_reward
57
58        plt.figure(figsize=(12, 8))
59
60        for lr, rewards in results.items():
61            smoothed_rewards = rolling_average(rewards, window_size=50)
62            plt.plot(smoothed_rewards, label=lr)
63
64        plt.title('Learning Curves for Different Learning Rates in PGB Algorithm')
65        plt.xlabel('Episodes')
66        plt.ylabel('Undiscounted Return')
67        plt.legend()
68        plt.grid(True)
69        plt.show()
70
```
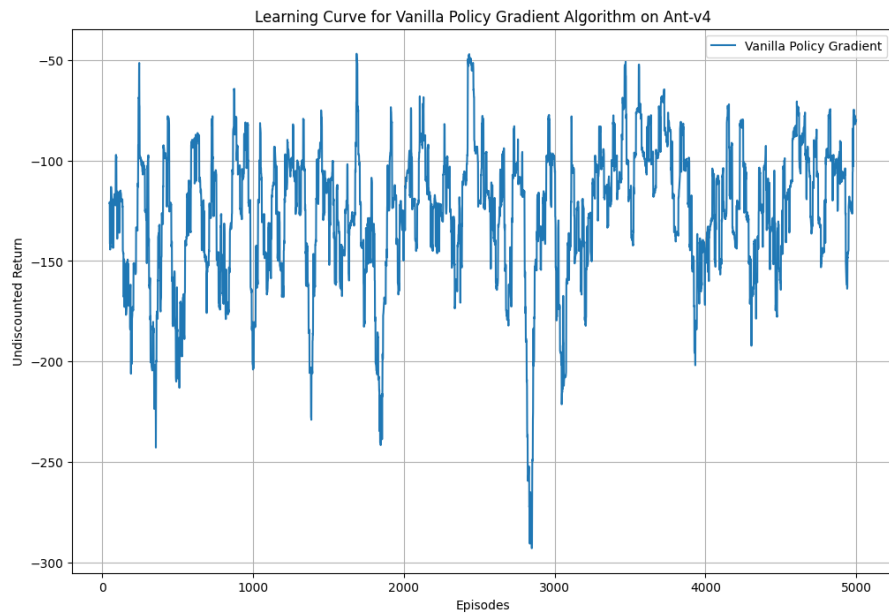
⌄ **PG algorithm output**

```
1 # Simulate command-line argument: replace 'pg' with 'pgb' to run policy gradient with baseline
2 sys.argv = ['hw2_F.py', '--algo', 'pg']
3
4 # Then call the main function
5 if __name__ == '__main__':
6     main()
```

```
    Episode:   0 Average Reward:  -38.3722569940707
    Episode:  1000 Average Reward:  -142.5412074282239
    Episode:  2000 Average Reward:  -133.20385312322185
    Episode:  3000 Average Reward:  -101.3894163135508
    Episode:  4000 Average Reward:  -153.6745736014923
```



Learning Curve for Vanilla Policy Gradient Algorithm on Ant-v4
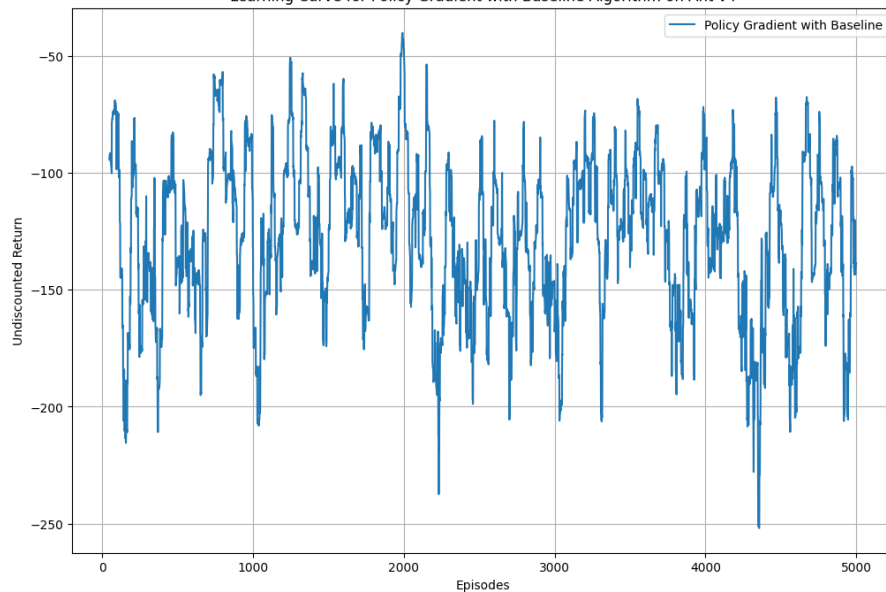
## PGB algorithm output

```
1 # Simulate command-line argument: replace 'pg' with 'pgb' to run policy gradient with baseline
2 sys.argv = ['hw2_F.py', '--algo', 'pgb']
3
4 # Then call the main function
5 if __name__ == '__main__':
6     main()
```

```
Episode:  0 Average Reward:  -59.20542941625349
Episode:  1000 Average Reward:  -114.72267109077255
Episode:  2000 Average Reward:  -92.31426428371581
Episode:  3000 Average Reward:  -152.1896180178551
Episode:  4000 Average Reward:  -103.55475125602347
```



Learning Curve for Policy Gradient with Baseline Algorithm on Ant-v4

```
Running experiment with learning rate: 0.0001
Episode:  0 Average Reward:  -74.51237787356047
Episode:  1000 Average Reward:  -116.63374839610125
Episode:  2000 Average Reward:  -141.0960428273539
Episode:  3000 Average Reward:  -100.1127510928168
Episode:  4000 Average Reward:  -167.94084120539
Running experiment with learning rate: 1e-05
Episode:  0 Average Reward:  -44.48651343751414
Episode:  1000 Average Reward:  -135.58542316701477
Episode:  2000 Average Reward:  -93.82546841978532
Episode:  3000 Average Reward:  -136.1494972191289
Episode:  4000 Average Reward:  -143.8363538322467
Running experiment with learning rate: 1e-06
Episode:  0 Average Reward:  -60.23545479633949
Episode:  1000 Average Reward:  -122.50935612246671
Episode:  2000 Average Reward:  -118.95391982287322
Episode:  3000 Average Reward:  -127.64069280982977
Episode:  4000 Average Reward:  -157.69421967281787
```



Learning Curves for Different Learning Rates in PGB Algorithm