

Assignment 3: Imitation Learning

Varun Shinde – 2024JRB2029

April 23, 2025

1 Algorithm Approach

We employ the DAgger (Dataset Aggregation) algorithm for imitation learning. DAgger is an iterative approach designed to address the limitations of behavior cloning. Instead of solely relying on expert trajectories (which the learner may not visit), the agent generates its own trajectories and queries the expert for corresponding actions. These expert-labeled states are then used to improve the learner's policy.

2 Environment Specifications

2.1 Observation Spaces

1. **Hopper-v4:** $(-\infty, \infty)$, shape: (11,), dtype: float64
 - **qpos (5 elements):** Position values of the robot's body parts.
 - **qvel (6 elements):** Velocity values (time derivatives of positions).
2. **Ant-v4:** $(-\infty, \infty)$, shape: (105,), dtype: float64
 - **qpos (13 elements):** Position values of the robot's body parts.
 - **qvel (14 elements):** Velocities of these body parts.
 - **cfrc_ext (78 elements):** External forces acting on each body part. Represented as a 13×6 array (force and torque components in 3D).

2.2 Action Spaces

1. **Hopper-v4:** $(-1.0, 1.0)$, shape: (3,), dtype: float64
 - A 3-dimensional vector representing torque applied to each joint.
2. **Ant-v4:** $(-\infty, \infty)$, shape: (8,), dtype: float64
 - An 8-dimensional vector representing motor commands to the ant's joints.

2.3 Evaluation:

We are making use of reward as a evaluation metric, while saving the best model. This is done since reward will be a appropriate metric to judge whether the current state of the learner policy is performing better than rest of iterations.

3 Pseudocode:

Algorithm 1 DAgger Training Loop

```
1: while not converged do
2:   Generate rollout using current learner policy  $\pi_\theta$  to collect state-action pairs
3:   Add generated trajectories to the replay buffer  $\mathcal{D}$ 
4:   Sample  $n$  trajectories from the replay buffer
5:   for each sampled state  $s_t$  do
6:     Query expert policy  $\pi^*$  to obtain expert action  $a_t^* = \pi^*(s_t)$ 
7:     Predict learner action  $a_t = \pi_\theta(s_t)$ 
8:     Compute MSE loss:  $\mathcal{L} = \frac{1}{n} \sum_{t=1}^n \|a_t - a_t^*\|^2$ 
9:   end for
10:  Update learner policy  $\pi_\theta$  using gradient descent on  $\mathcal{L}$ 
11:  Evaluate current policy  $\pi_\theta$  by generating evaluation trajectories
12:  Compute average reward  $R_{\text{avg}}$ 
13:  if  $R_{\text{avg}} > R_{\text{best}}$  then
14:    Save model checkpoint
15:     $R_{\text{best}} \leftarrow R_{\text{avg}}$ 
16:  end if
17: end while
```

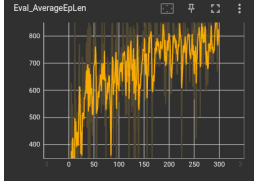
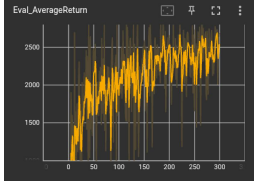
Visualization	Inference
	This graph shows the average episode length. The increasing trend indicates that the agent is learning to sustain longer episodes, reflecting improved stability and performance.
	This plot represents the average return (cumulative reward). The upward trend suggests that the learner is successfully mimicking the expert and achieving higher rewards, validating the effectiveness of the imitation learning approach.

Table 1: Evaluation Metrics During DAgger Training on Hopper-v4 Environment

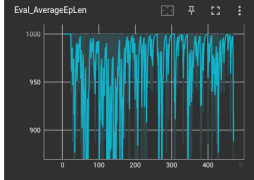

Visualization	Inference
	This graph shows the average episode length. The increasing trend indicates that the agent is learning to sustain longer episodes, reflecting improved stability and performance.
	This plot represents the average return (cumulative reward). The upward trend suggests that the learner is successfully mimicking the expert and achieving higher rewards, validating the effectiveness of the imitation learning approach.

Table 2: Evaluation Metrics During DAgger Training on Ant-v4 Environment