

Supervisor: Anna Lukina

#### 1. Introduction

Imitation learning trains a policy which attempts to imitate the expert. When using imitation learning in critical decision-making processes, it is necessary for these policies to be interpretable.

We can train an interpretable surrogate model using imitation learning from the expert black boxes to obtain the required policies.

How do decision trees, produced by the AggreVaTe[1] algorithm, compare to a **Behavioral Cloning** baseline, and other imitation learning algorithms in terms of interpretability and performance?

The other imitation learning algorithms are: GAIL[2], Viper[3] and DAgger[4]

#### 2. Methodology

- AggreVaTe is an imitation learning algorithm that uses cost-togo to obtain better long-term data
- The expert calculates, for each possible action, the cost to get to the finish when performing that action



Figure 1: AggreVaTe algorithm learning process

- AggreVaTe requires a high number of training iterations or data collection rollouts due to obtaining only one datapoint per rollout
- The AggreVaTe algorithm has been modified to train decision tree policies to make it possible to compare on interpretability

# Interpretability and performance comparisons of decision tree surrogate models produced by AggreVaTe

### Setup

- 3 Open Al Gym Environments
- Fixed parameters
  - max depth of tree
  - number of iterations
  - rollouts per iteration



### 4.Conclusion

In terms of interpretability, AggreVaTe performs equal or better than all other imitation learning algorithm.

In terms of performance, AggreVaTe performs only slightly worse than GAIL, Viper and DAgger, however it performs better than Behavioral Cloning.

This is consistent with expectations since fewer data points lead to more possible failing paths that have not yet been explored, but it also leads to a decision tree less prone to overfitting.

3. Experiments

#### Policy • A decision tree model (figure 2) • Shows all decisions made • Usable to find an explanation for difference in performance BCMountain Car Average Reward -120.8Standard deviation 4.0Nodes Table 1: Results for mountain car with $\mathbf{BC}$ Cartpole PoleAngle ≤ -0.025 Average Reward 207.276.6% Standard deviation 44.2[0.153, 0.847]Nodes Right Table 2: Results for cartpole with depth BCAcrobot Average Reward -94.567.8% Standard deviation 37.3[0.092, 0.908]Nodes Right Table 3: Results for acrobot with depth 2

[1] Stephane Ross and J Andrew Bagnell. Reinforcement and [4] Stephane Ross, Geoffrey Gordon, and Drew imitation learning via interactive. no-regret learning. Bagnell. A reduction of imitation learning and structured prediction to no-regret online learning. arXiv preprint arXiv:1406.5979, 2014. [2] Jonathan Ho and Stefano Ermon. Generative adversarial In Proceedings of the fourteenth international imitation learning. Advances in neural information processing conference on artificial intelligence and statistics, systems, 29:4565–4573, 2016. pages 627–635. JMLR Workshop and Conference [3] Osbert Bastani, Yewen Pu, and Armando Solar-Lezama. Proceedings, 2011. Verifiable reinforcement learning via policy extraction. arXiv preprint arXiv:1805.08328, 2018.



## Results

Comparison using the metrics:

- Average reward
- Standard deviation
- Number of nodes

The outcomes can be found in table 1, 2 and 3, where the best are in bold.

DAgger	GAIl	VIPER	AggreVaTe
-120.3	-119.2	-119.6	-120.4
<b>3.6</b>	3.8	3.7	4.5
3	3	3	3
depth 1			
DAgger	GAIl	Viper	AggreVaTe
499.8	498.2	<b>500</b>	488.4
0.7	14.12	0	22.0
7	11	15(9)	5
h 3			
DAgger	GAIl	Viper	AggreVaTe
-83.5	-83.0	-84.4	-85.5
<b>10.8</b>	14.9	21.0	24.5
<b>5</b>	7	7	3

#### 5. Limitations

Experiments were on simple environments Experiments were performed with different experts The high number of rollouts required could have caused overfitting of the other algorithms.

#### References