# Interpretability of **Surrogate Models Produced by Viper**

Otto Kaaij Supervisor: Anna Lukina

# Background

- There is a need for interpretable AI
- One possible approach:

Use imitation learning to extract interpretable surrogate models from black box oracles







**Black box** oracle  $\pi^*$ 

Imitation learning

Interpretable decision tree  $\pi$ 

Left

 $x \le 0.391$ 

# Viper

Viper<sup>[1]</sup> is an imitation learning algorithm that prioritizes accuracy on critical states: states where making the wrong choice has the highest cost

# Goals

- Evaluate decision trees created by Viper on performance and interpretability
- Compare Viper to other imitation learning algorithms (Behavioral cloning (BC), Gail<sup>[2]</sup> and AggreVate<sup>[3]</sup>)

# **Results**



Right

X

Å

 $\pi$  reward:

-111.2±2.7

False

 $\theta'_{1} \leq -0.23$ 

-1 T

+1 T

 $\theta'_2 \le 0.155$ 

+1 T

trees

# **Future Work**

- Evaluate more complex environments • Evaluate other types of surrogate models • Improve Viper by training decision trees more effectively and by cross-validating on interpretability • Unify oracles for better comparisons

# References

2016



Right Left Right Nothing Interpretation Accelerate in current direction to build  $\pi^*$  reward:  $\textcircled{Q}_{3}$ momentum. Reverse just -112.1±1.8 before standing still.

 $x \le -0.37$ 



# **TUDelft**

## **Algorithm Comparisons**

- Viper, Gail and AggreVaTe all outperform baseline behavioral cloning
- Viper is the only algorithm that matches or
  - improves on oracle performance on all three environments
- All algorithms produce similarly interpretable results, though Viper tends to produce smaller

# Conclusion

• Viper produces high performance, interpretable decision trees for simple environments • Performance is dependent on oracle quality, which cannot be expressed only in oracle reward • We can use this interpretability to understand the policies and, in some cases, improve them

- [1] Osbert Bastani, Yewen Pu, and ArmandoSolar-Lezama. "Verifiable reinforcement learning via policy extraction". In:arXiv preprintarXiv:1805.08328(2018) [2] Jonathan Ho and Stefano Ermon. "Generative Adversarial Imitation Learning". In: Advances in Neural Information Processing Systems. Ed. by D.vLee et al. Vol. 29. Curran Associates, Inc.,
- [3] Stephane Ross and J Andrew Bagnell. "Reinforcement and imitation learning via interactive no-regret learning". In: arXiv preprint arXiv:1406.5979 (2014). Icons under CC by Milinda Courey, DinosoftLab and M. Oki Orlando from NounProject.com