

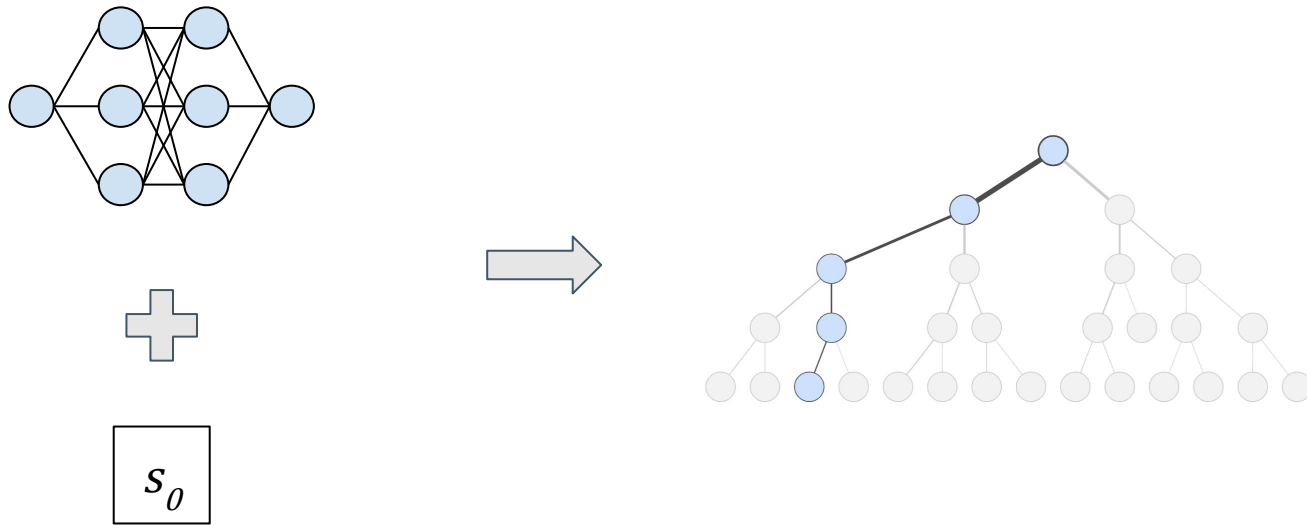
Interpretable Local Tree Surrogate Policies

AAAI SAFEAI WORKSHOP 2022

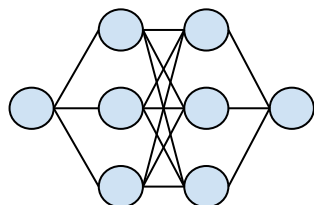
John Mern, Sidhart Krishnan, Anil Yildiz, Kyle Hatch,
Mykel J Kochenderfer

SISL
Stanford Intelligent
Systems Laboratory

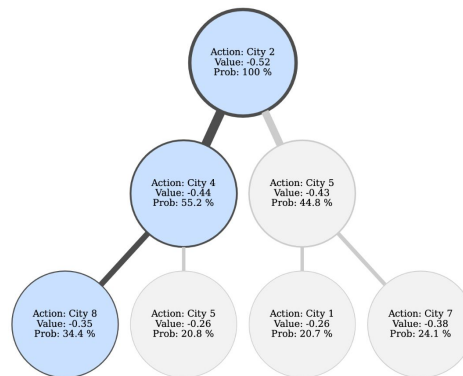
Trees as Interpretable Surrogates



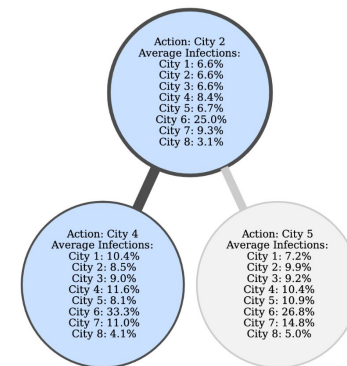
Trees as Interpretable Surrogates



S_0

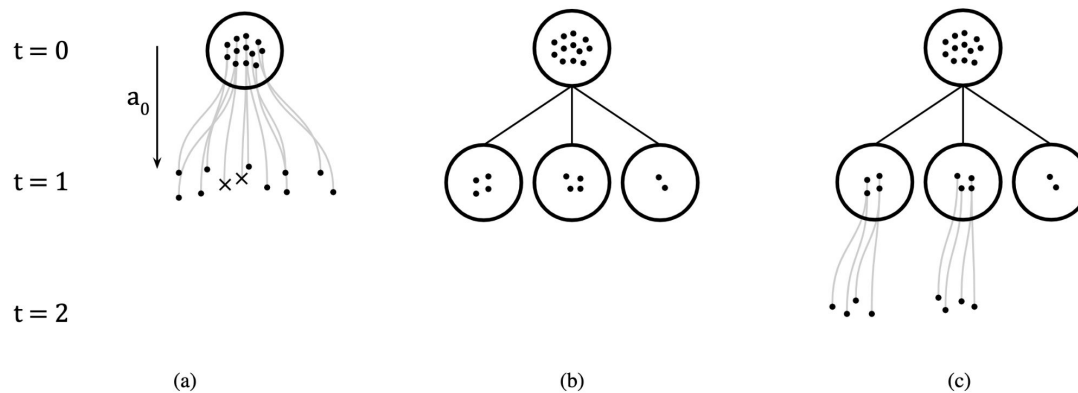


(a)



(b)

Monte Carlo Clustering



$$\delta = \left\| \hat{Q}(s, a) - \max_{a' \in \mathcal{A}} \hat{Q}(s, a') \right\|$$

SISL

Stanford Intelligent
Systems Laboratory