





Explanation-based Finetuning Makes Models More Robust to Spurious Cues

Josh Magnus Ludan, Yixuan Meng*, Tai Nguyen*, Saurabh Shah*, Qing Lyu, Marianna Apidianaki, Chris Callison-Burch



60







Main Results



Analysis: Cue Pervasiveness

 On Embedding cue (cluster data in 2 parts unsupervised), explanation-based finetuning also scales well with the strength of the cue.



Additional Takeaways

 Finetune with intentionally FALSE explanations still mitigates the correlation better that finetune *without* explanations. • Our method suggests a strong synergy between interpretability and robustness.