Automated Detection of Visual Attribute Reliance with a Self-Reflective Agent

Christy Li, Josep Lopez Camuñas, Jake Touchet, Jacob Andreas,

Agata Lapedriza, Antonio Torralba, Tamar Rott Shaham





On Machine Learning

Motivation

- A well-known failure mode of vision models performing image recognition is over-reliance on spurious visual attributes, resulting in brittle predictions that fail to generalize.
- Existing methods for discovering the visual attributes behind model predictions often rely on manual inspection along a fixed set of predefined features, making them difficult to scale.
- We propose a Self-reflective Automated Interpretability Agent (SAIA) that produces natural-language descriptions of spurious attribute dependencies of visual models.

Attribute Reliance Detection Benchmark

Bird detector that relies on the presence of beach background



Self-Reflection



Taxonomy of Visual Attribute Reliances \mathscr{A}_i

Object attributes: color, material

Context attributes: image background, object state/configuration

Demographic attributes: gender, age of human subjects

Counterfactual demographic attributes: gender, age of human subjects

130 unique object recognition models, 18 types of visual attribute reliance

SAIA Framework



show lower classification score than expected, let's continue with the experiments. Highest classification scores were for direct child interaction, rather than educational settings in general. Educational setups without human presence scored only moderately. Isolated bears in various settings and environmental conditions didn't significantly impact scores.



Evaluation

How does predictiveness score change with increasing rounds of self-reflection?



Applications to SOTA Models



> Test: corporate settings



The model relies on corporate settings.

> Looking at images from self-evaluation and generating new hypotheses.

> Test: 45-degree views of suits



> Test: male vs. female



The model exhibits a strong gender bias, relying on a male wearing the suit.



Predicted low-scoring images

Some images with corporate settings got low confidence scores, while some without corporate settings got high scores. Let's run another experiment.

> Let's evaluate the "gender bias" conclusion.

Predicted high-scoring images





The conclusion matches the model behavior: The model relies on male presenting figures to detect the suit.



What are the visual attributes that CLIP relies on to recognize the concept teacher?



Let's conduct some experiments...

CLIP relies on authority signals in conventional classroom settings rather than actual pedagogical engagement.

Predicted high-scoring images





Predicted low-scoring images



What are the visual attributes that YOLO relies on to detect pedestrians?



Let's conduct some experiments...

YOLO relies on side-view walking poses in urban crossing contexts.

Predicted high-scoring images





Predicted low-scoring images

We demonstrate that SAIA can identify feature reliances in state-of-the-art models including **CLIP's** vision encoder and the **YOLOv8** object detector.