

21

How to Read the "Mind" of a Neural Network

Thomas Walton, Advisor: Mohammad R. Hasan



Introduction & Motivation

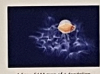
- Deep Learning (DL) is a powerful tool for solving computer vision problems, but has shortcomings that are not well understood
- Algorithmic bias plagues the decision-making process of DL models, leading to nongeneralizable results and limited capability
- A Convolutional Neural Network (CNN) is a DL model used to learn the semantic identity (category) of data
- CNNs are intelligent (they can generate knowledge learned from data onto novel, unseen data)
- Algorithmic bias limits the capability of CNNs to learn



A CNN analyzing an image of a cat

Objective

- Use Class Activation Maps (CAMs) to read the "mind" of a CNN
- Create ScoreCAM maps to make these CAMs human readable
- Leverage this analysis to discover algorithmic bias



A ScoreCAM map of a daisy flower shows where a model was learning to decide the class of this image

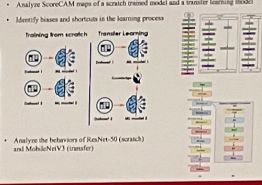
Scientific Research Questions

- SRQ1: What are the reasons for a vision model to fail in its predictions?
- SRQ2: When a vision model identifies an object accurately, does it necessarily mean that it recognizes the object in the image? If not, then why?



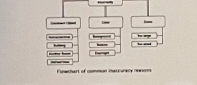
This decision-making process is filtered through layers of neurons, each contributing to the final prediction

Methodology



Results

SRQ1: Both ResNet-50 and MobileNetV3 had similar drawbacks



CNNs can make mistakes by recognizing features, rather than by recognizing what they're trying to recognize in the scene

Conclusions

- On their own, SRQ1 and SRQ2 tell only part of the story of why algorithmic bias occurs
- When failure (from SRQ1) are combined with suspicious from SRQ2, bias is exposed
- Human error
- Adding unrepresentative objects lead to misclassification on images where the model was once confident
- By adding objects that lead to failure onto images where the vision model was confident, weaknesses in CNNs become apparent
- The lack of adaptability from this vision model implies that while it may be good at classifying flowers, it is not truly learning intelligent representations



SRQ2: While a model may seem to have learned intelligence, further analysis shows that this may not be true



Despite being correct in their predictions, both models used obscure or other objects to draw conclusions



While the CNN correctly identifies the image, when shown an, trouble arises