

Artificial intelligence must know when to ask for human help

March 7 2019, by Sarah Scheffler, Adam D. Smith And Ran Canetti



Credit: AI-generated image (disclaimer)

Artificial intelligence systems are powerful tools for businesses and governments to process data and respond to changing situations, whether <u>on the stock market</u> or <u>on a battlefield</u>. But there are still some things AI isn't ready for.



We are <u>scholars of computer science working to understand</u> and improve the ways in which algorithms interact with society. AI systems perform best when the goal is clear and there is high-quality data, like when they are asked to distinguish between different faces after learning from many pictures of correctly identified people.

Sometimes AI systems do so well that users and observers are surprised at <u>how perceptive</u> the technology is. However, sometimes success is <u>difficult to measure</u> or <u>defined incorrectly</u>, or the training data <u>does not</u> <u>match the task at hand</u>. In these cases, AI algorithms tend to fail in <u>unpredictable and spectacular ways</u>, though it's <u>not always immediately</u> <u>obvious</u> that something has even gone wrong. As a result, it's important to be wary of the hype and excitement about what AI can do, and not assume the solution it finds is always correct.

When algorithms are at work, there should be a human safety net to prevent harming people. Our research demonstrated that in some situations algorithms can recognize problems in how they're operating, and <u>ask for human help</u>. Specifically, we show, asking for human help can help alleviate algorithmic bias in some settings.

How sure is the algorithm?

Artificial intelligence systems are being used in <u>criminal sentencing</u>, <u>facial-based personality profiling</u>, <u>resume screening</u>, <u>health care</u> <u>enrollment</u> and other difficult tasks where people's lives and well-being are at stake. U.S. government agencies are beginning to ramp up their exploration and use of AI systems, in response to a recent <u>executive</u> <u>order from President Donald Trump</u>.

It's important to remember, though, that AI can cement misconceptions in how a task is addressed, or magnify existing inequalities. This can happen even when no one told the <u>algorithm</u> explicitly to treat anyone



differently.

For instance, many companies have algorithms that try to determine features about a person by their face – say to guess their gender. The systems developed by U.S. companies tend to do significantly <u>better at categorizing white men</u> than they do women and darker-skinned people; they do worst at dark-skinned women. Systems developed in China, however, tend to <u>do worse on white faces</u>.

The difference is not because one group has faces that are easier to classify than others. Rather, both algorithms are typically trained on a large collection of data that's not as diverse as the overall human population. If the data set is dominated by a particular type of face – white men in the U.S., and Chinese faces in China – then the algorithm will probably do better at analyzing those faces than others.





Credit: Unsplash/CC0 Public Domain

No matter how the difference arises, the result is that algorithms can be biased by being more accurate on one group than on another.

Keeping a human eye on AI

For high-stakes situations, the algorithm's confidence in its own result – its estimation of how likely it is that the system came up with the right answer – is just as important as the result itself. The people who receive the output from algorithms need to know how seriously to take the results, rather than assuming that it's correct because it involved a computer.

Only recently have researchers begun to develop ways to identify, much less attempt to fix, <u>inequalities in algorithms and data</u>. Algorithms can be programmed to recognize their own shortcomings – and follow that recognition with a <u>request for a person to assist with the task</u>.

Many types of AI algorithms already calculate an internal <u>confidence</u> <u>level</u> – a prediction of how well it did at analyzing a particular piece of input. In facial analysis, many AI algorithms <u>have lower confidence</u> on darker faces and female faces than for white male faces. <u>It's unclear</u> how much this has been taken into account by law enforcement for highstakes uses of these algorithms.

The goal is for the AI itself to locate the areas where it is not reaching the same accuracy for different groups. On these inputs, the AI can defer its decision to a human moderator. This technique is especially



well-suited for context-heavy tasks like content moderation.

Human content moderators <u>cannot keep up with</u> the flood of images being posted on social media sites. But AI content moderation is famous for failing to take into account the context behind a post – misidentifying discussions of sexual orientation as <u>explicit content</u>, or identifying the Declaration of Independence as <u>hate speech</u>. This can end up inaccurately censoring one <u>demographic</u> or <u>political</u> group over another.

To get the best of both worlds, <u>our research</u> suggests scoring all content in an automated fashion, using the same AI methods already common today. Then our approach uses newly proposed techniques to automatically locate potential inequalities in the accuracy of the algorithm on different protected groups of people, and to hand over the decisions about certain individuals to a human. As a result, the algorithm can be completely unbiased about those people on which it actually decides. And humans decide on those individuals where algorithmic decision would have inevitably created bias.

This approach does not eliminate bias: It just "concentrates" the potential for bias on a smaller set of decisions, which are then handled by people, using human common sense. The AI can still perform the bulk of the decision-making work.

This is a demonstration of a situation where an AI algorithm working together with a human can reap the benefits and efficiency of the AI's good decisions, without being locked into its bad ones. Humans will then have more time to work on the fuzzy, difficult decisions that are critical to ensuring fairness and equity.

This article is republished from <u>The Conversation</u> under a Creative Commons license. Read the <u>original article</u>.



Provided by The Conversation

Citation: Artificial intelligence must know when to ask for human help (2019, March 7) retrieved 28 June 2024 from <u>https://phys.org/news/2019-03-artificial-intelligence-human.html</u>

This document is subject to copyright. Apart from any fair dealing for the purpose of private study or research, no part may be reproduced without the written permission. The content is provided for information purposes only.