# Developing a moral compass from human texts

February 7 2019



Can machines develop a moral compass? Credit: Patrick Bal

Artificial Intelligence (AI) translates documents, suggests treatments for patients, makes purchasing decisions and optimises workflows. But where is its moral compass? A study by the Centre for Cognitive Science at TU Darmstadt shows that AI machines can indeed learn a moral

compass from humans. The results of the study have been presented at this year's ACM/AAAI Conference on AI, Ethics, and Society (AIES).

AI has an increasing impact on our society. From self-driving cars on [public roads](link), to self-optimising industrial production systems, to health care – AI machines handle increasingly complex human tasks in increasingly autonomous ways. And in the future, autonomous machines will appear in more and more areas of our daily lives. Inevitably, they will be confronted with difficult decisions. An autonomous robot must know that it should not kill people, but that it is okay to kill time. The robot needs to know that it should rather toast a slice of bread than a hamster. In other words: AI needs a human-like moral compass. But can AI actually learn such a compass from humans?

Researchers from Princeton (USA) and Bath (UK) had pointed out (*Science*, 2017) the danger that AI, when applied without care, can learn [word associations](link) from written texts and that these associations mirror those learned by humans. For example, the AI interpreted male names that are more common in the Afro-American community as rather unpleasant and names preferred by Caucasians as pleasant. It also linked female names more to art and male names more to technology. For this, huge collections of written texts from the internet were fed into a [neural network](link) to learn vector representations of words – coordinates, i.e. words get translated into points in a high-dimensional space. The [semantic similarity](link) of two words is then computed as the distance between their coordinates, the so-called [word embeddings](link), and complex semantic relations can be computed and described by simple arithmetic. This applies not only to the harmless example "king – man + woman = queen" but also to the discriminating "man – technology + art = woman".

## Machines can reflect our values

Now, a team led by professors Kristian Kersting and Constantin

Rothkopf at the Centre for Cognitive Science of the TU Darmstadt has successfully demonstrated that machine learning can also extract deontological, ethical reasoning about "right" and "wrong" conduct from written text. To this end, the scientists created a template list of prompts and responses, which include questions such as "Should I kill people?", "Should I murder people?", etc. with answer templates of "Yes, I should" or "No, I should not." By processing a large body of human texts the AI system then developed a human-like moral compass. The moral orientation of the machine is calculated via embedding of the questions and answers. More precisely, the machine's bias is the difference of distances to the positive response ("Yes, I should") and to the negative response ("No, I should not"). For a given moral choice overall, the model's bias score is the sum of the bias scores for all question/answer templates with that choice. In the experiments, the system learned that you should not lie. It is also better to love your parents than to rob a bank. And yes, you should not kill people, but it is fine to kill time. You should also put a slice of bread in the toaster rather than a hamster.

The study provides an important insight to a fundamental question in AI: Can machines develop a [moral compass](#)? And if so, how can we effectively "teach" machines our morale? The results show that machines can reflect our values. They can adopt human-like prejudices, indeed, but they can also adopt our moral choices by "observing" humans. In general, embeddings of questions and answers can be seen as a type of microscope that allow one to study the moral values of text collections as well as the development of moral values in our society.

The results from the study provide several avenues for future work, in particular when incorporating modules constructed via [machine learning](#) into decision-making systems.

**More information:** The Moral Choice Machine: Semantics Derived Automatically from Language Corpora Contain Human-like Moral

Aylin Caliskan et al. Semantics derived automatically from language corpora contain human-like biases, *Science* (2017). DOI: 10.1126/science.aal4230

Provided by Technische Universitat Darmstadt